



**Barcelona
Supercomputing
Center**
Centro Nacional de Supercomputación



esiwace
CENTRE OF EXCELLENCE IN SIMULATION OF WEATHER
AND CLIMATE IN EUROPE

Improving Ocean Model Computational Performance by using Mixed-Precision Approaches

Oriol Tintó, Mario Acosta, Miguel Castrillo,
Kim Serradell, Francisco J. Doblas-Reyes

Computational Earth Sciences

PASC19 - High-Resolution Weather and Climate
Simulations

14/06/2019

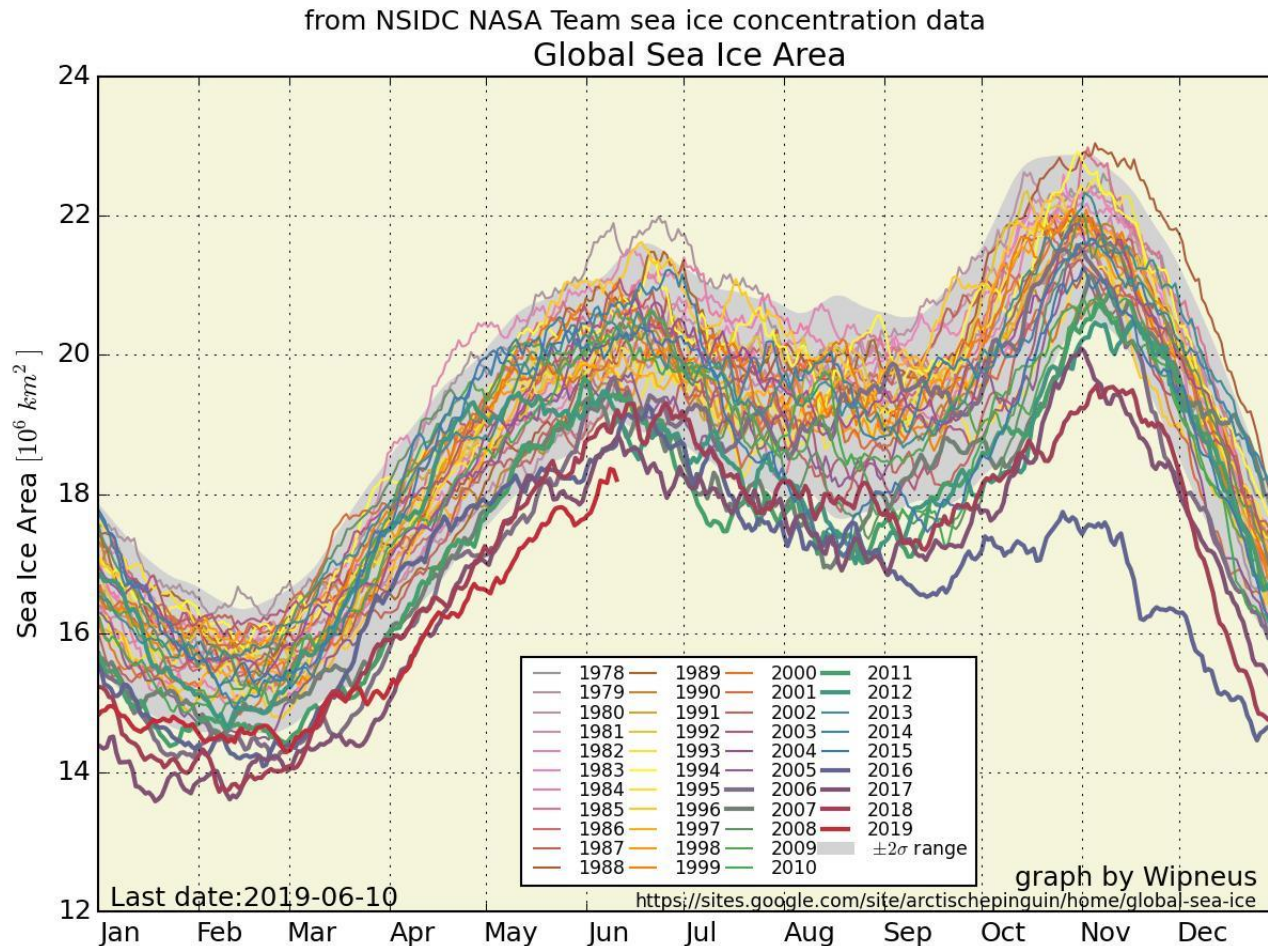
Outline

- Why?
- How?



esiwace
CENTRE OF EXCELLENCE IN SIMULATION OF WEATHER
AND CLIMATE IN EUROPE

Why models?



Why mixed precision?

0 nz

2,869,753,461 nz

96 processes

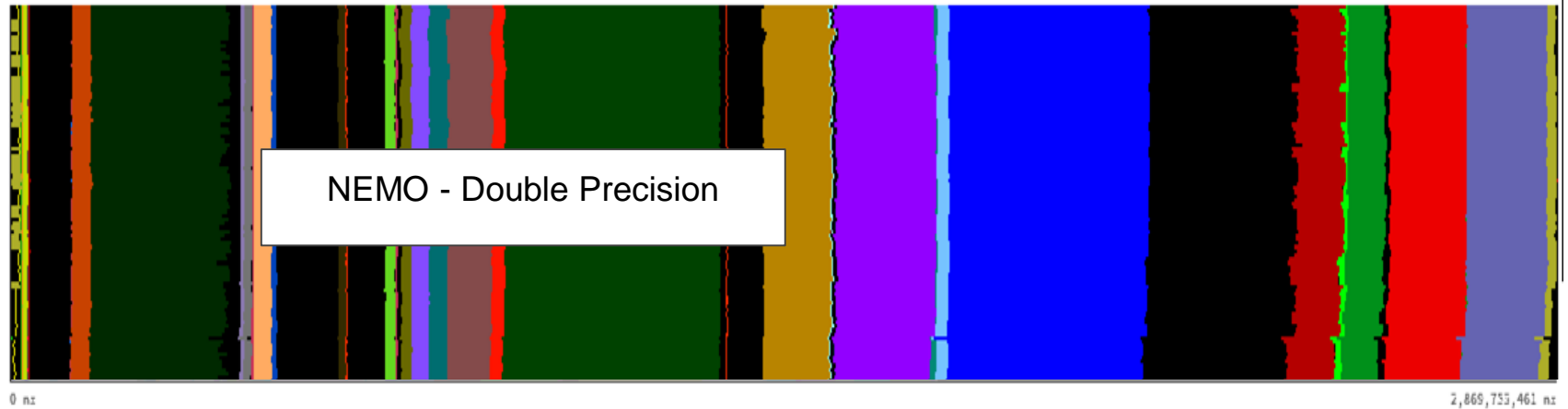
NEMO - Double Precision

0 ns

2,869,753,461 ns

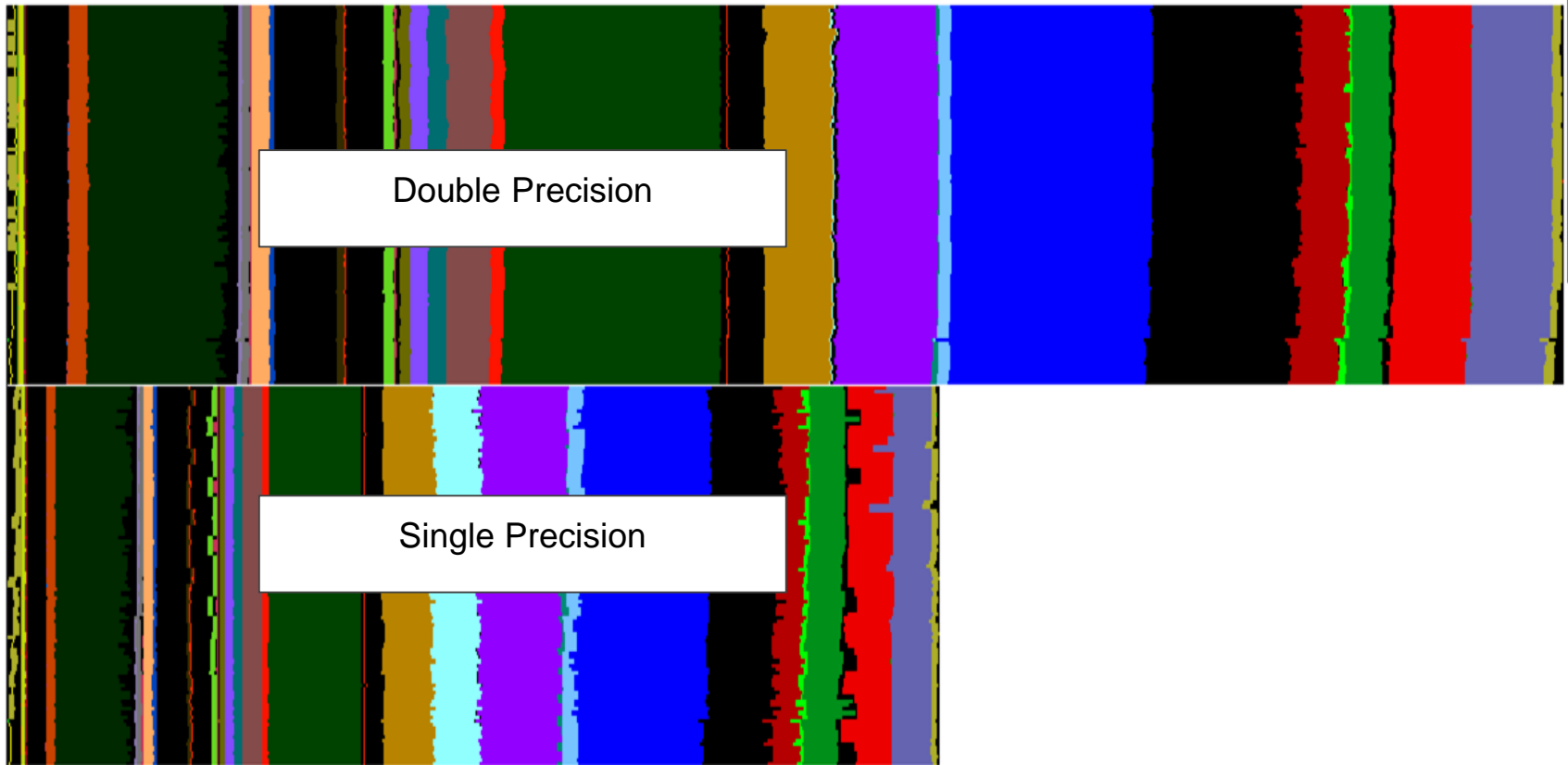
Time

96 processes



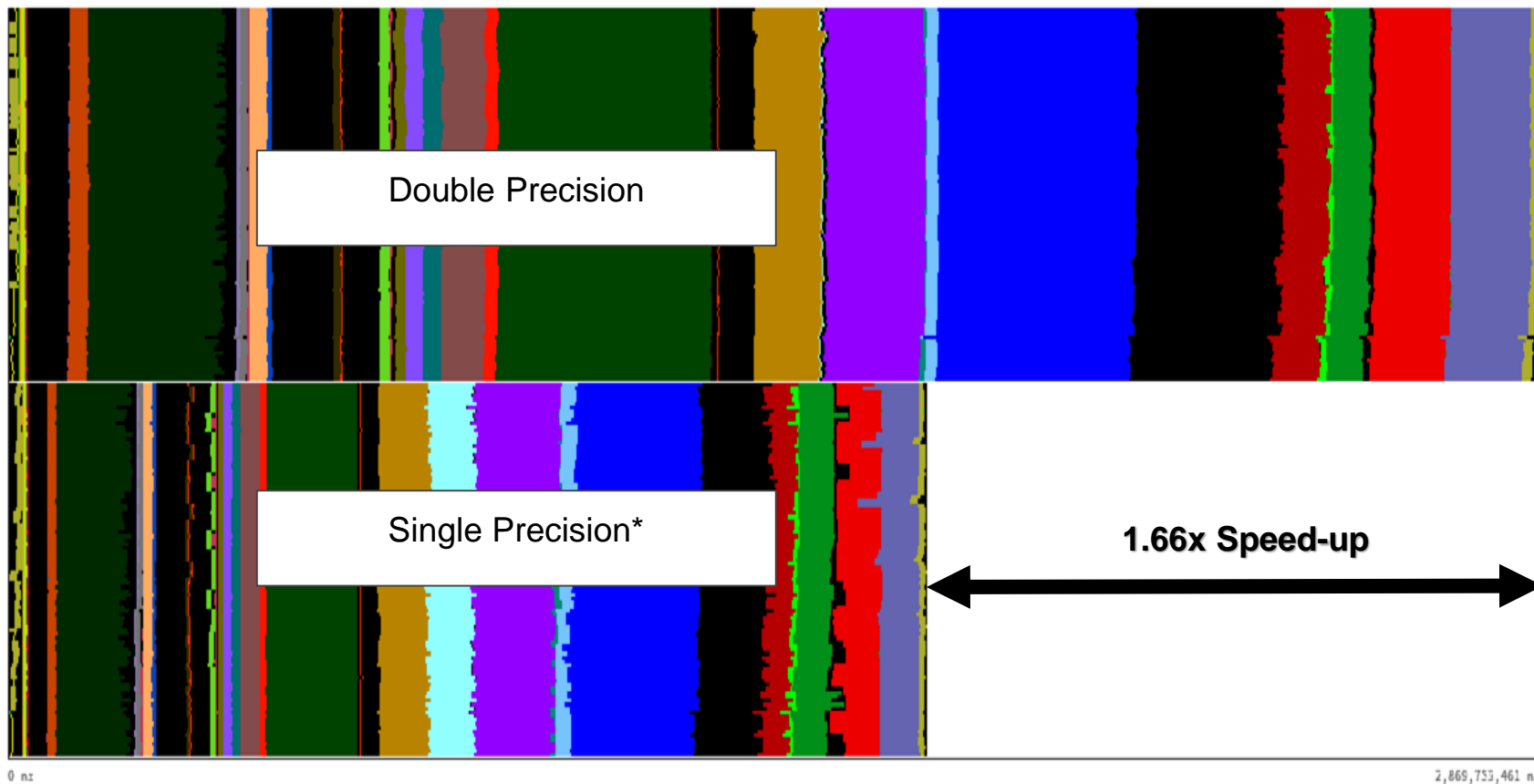
Time

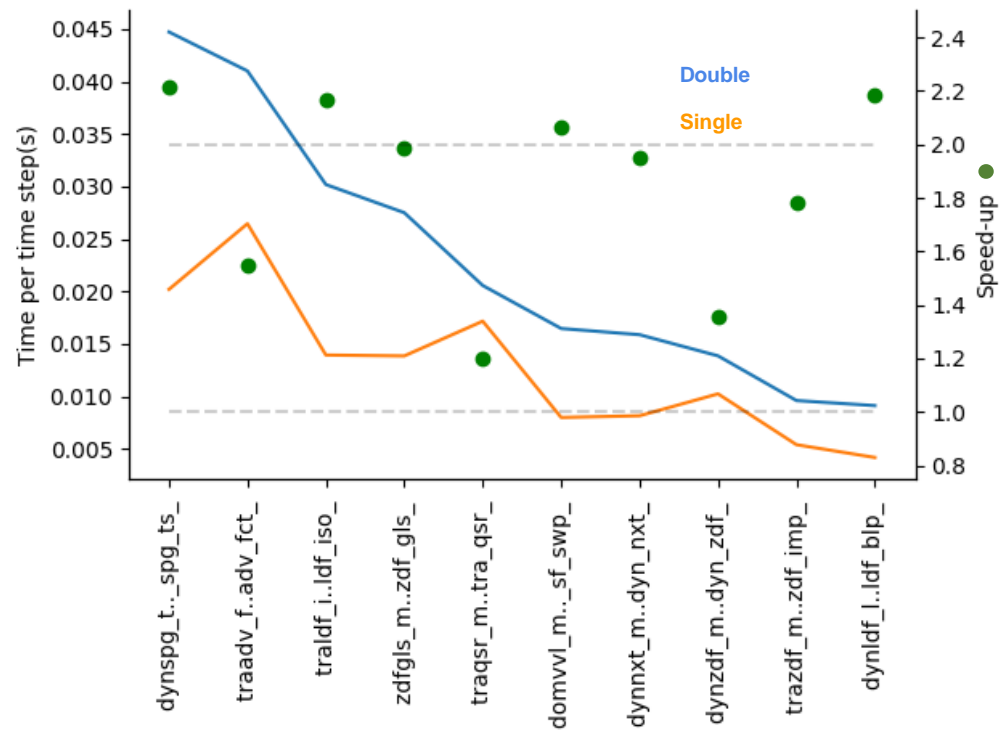
- The trace corresponds to a single time-step
- NEMO 4.0 Ocean-only
- ORCA025 -> global $\frac{1}{4}^\circ \sim 27\text{km}$ at equator
- Each color represents a different routine

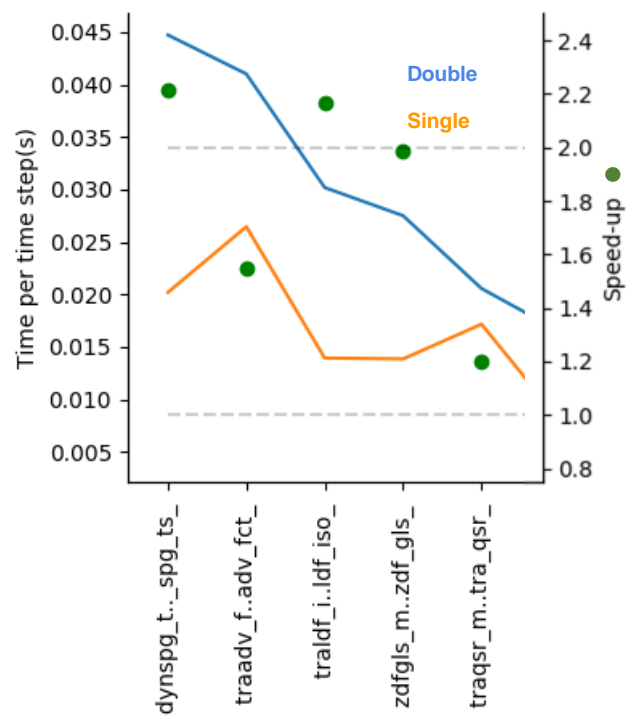


0 nz

2,869,753,461 nz







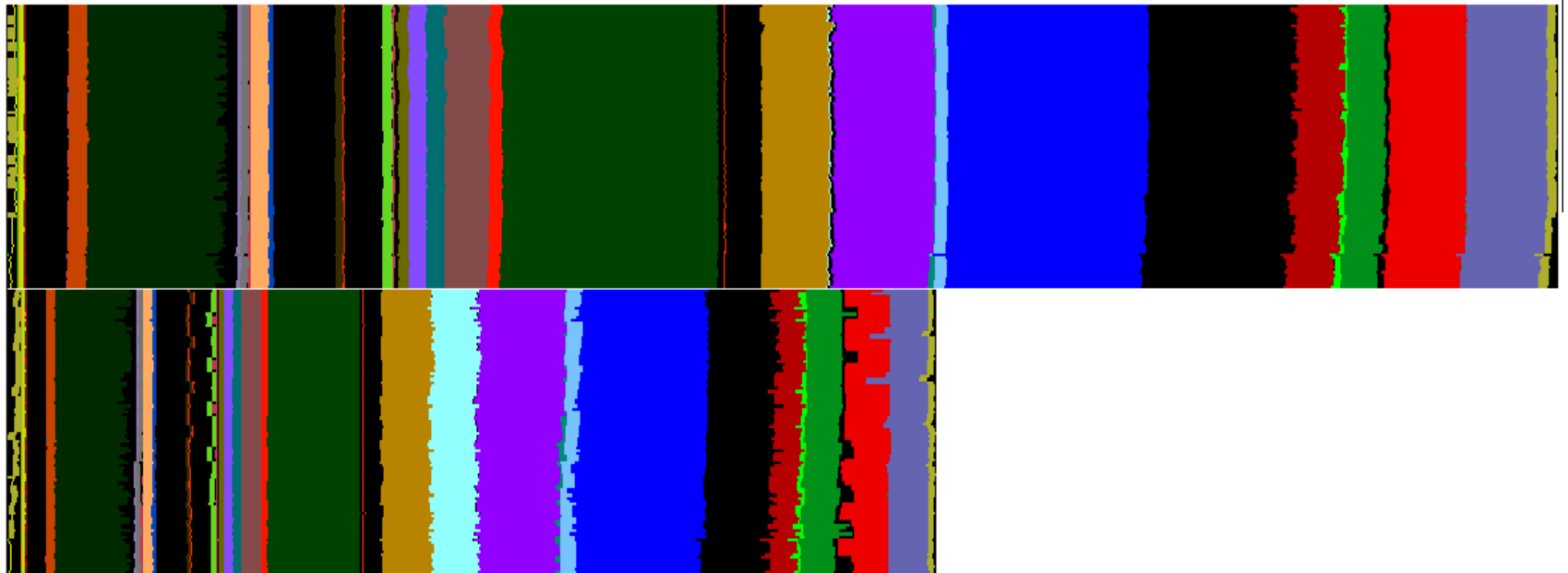
dynspg

traadv

traldf

zdfgls

traqsr



0 m2

2,868,753,461 m2

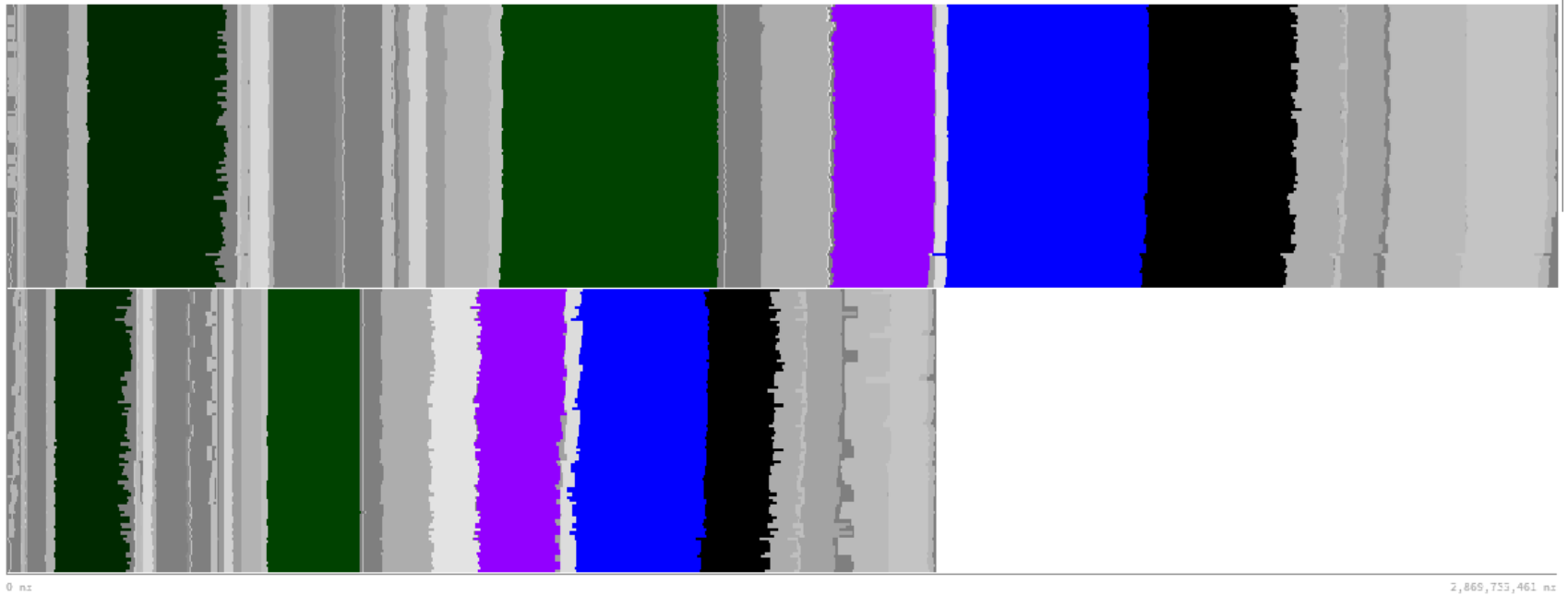
dynspg

traadv

traldf

zdfgls

traqsr



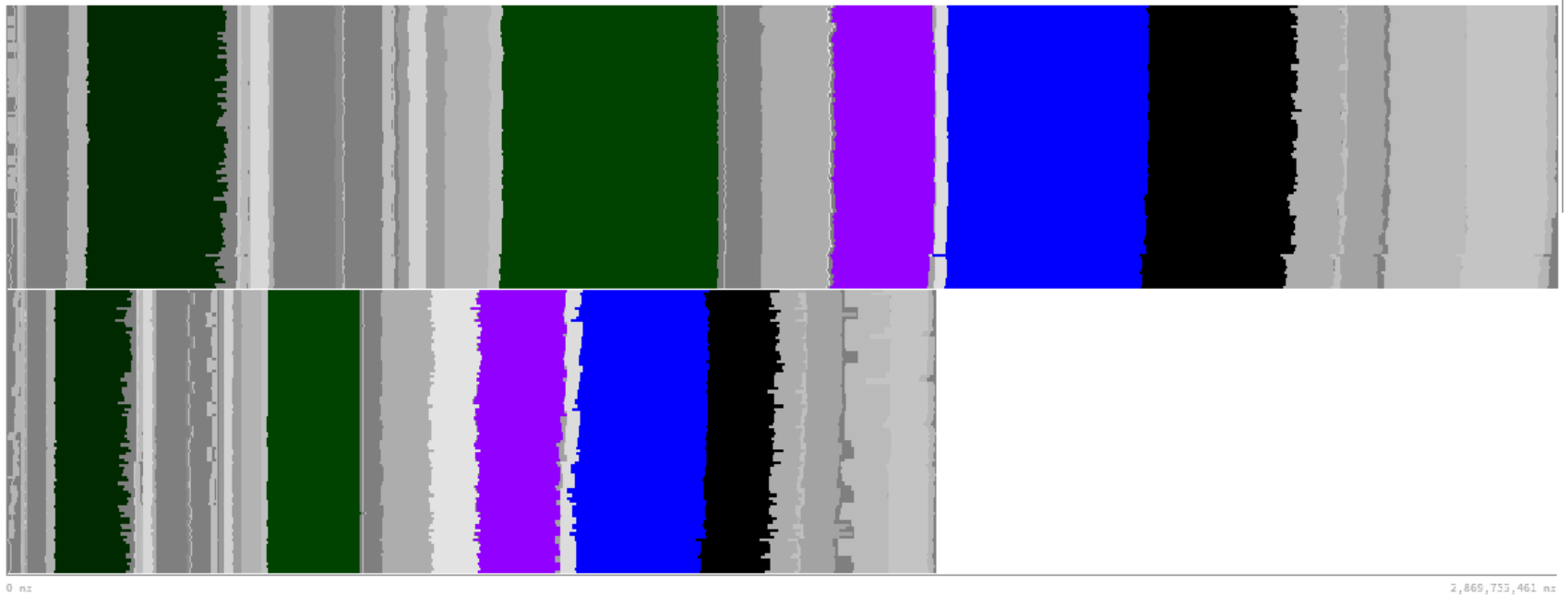
dynspg

traadv

traldf

zdfgls

traqsr



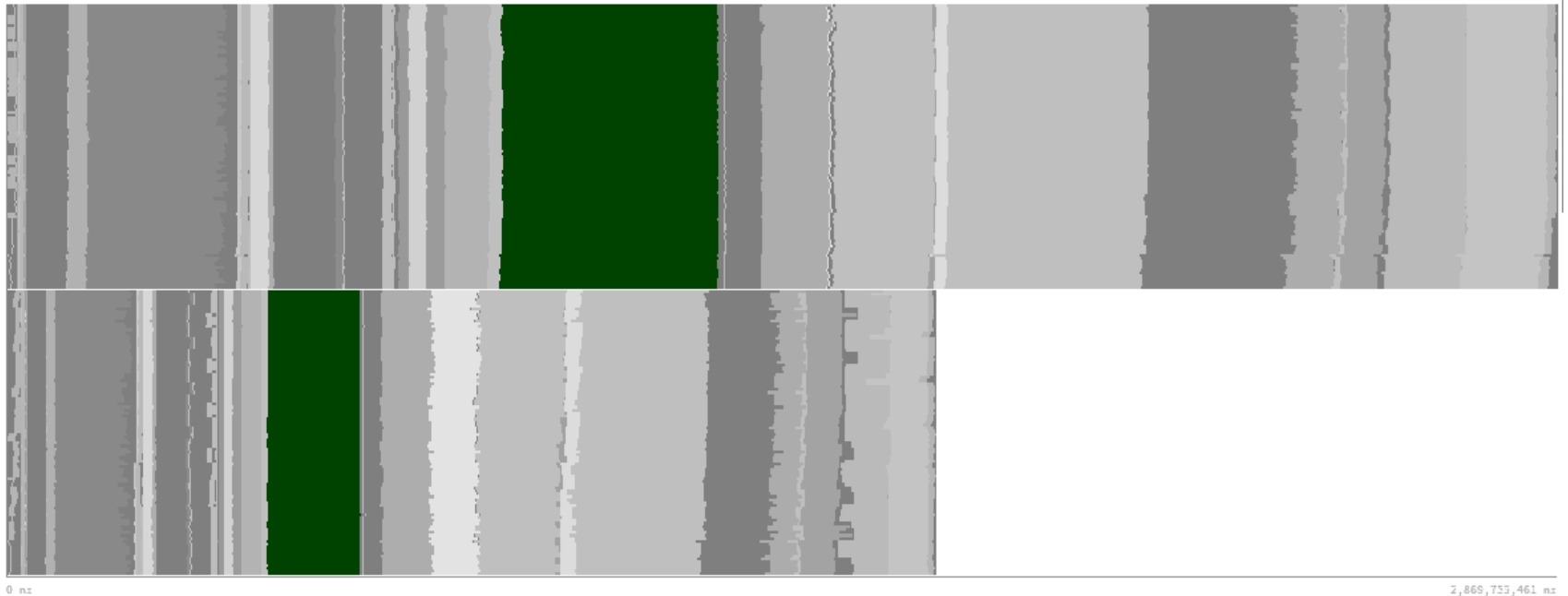
dynspg

traadv

traldf

zdfgls

traqsr



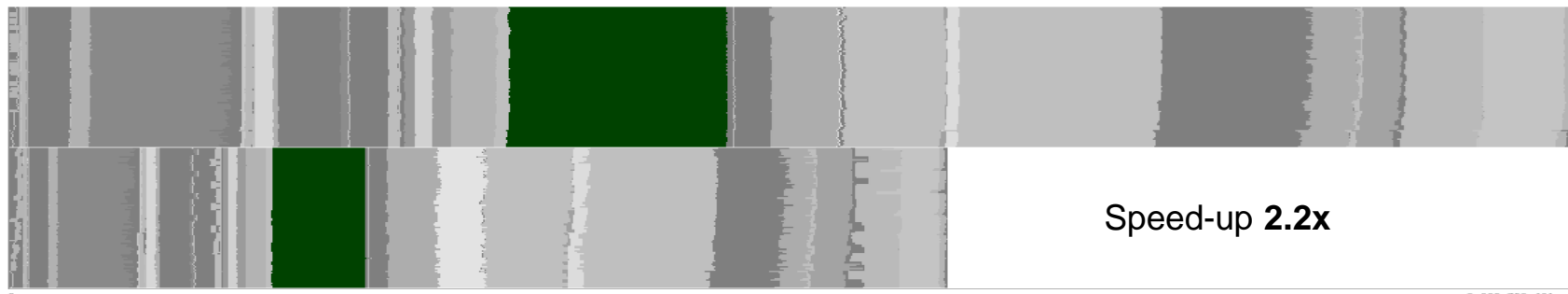
dynspg

traadv

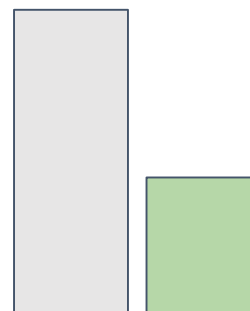
traldf

zdfgls

traqsr



Iteration time



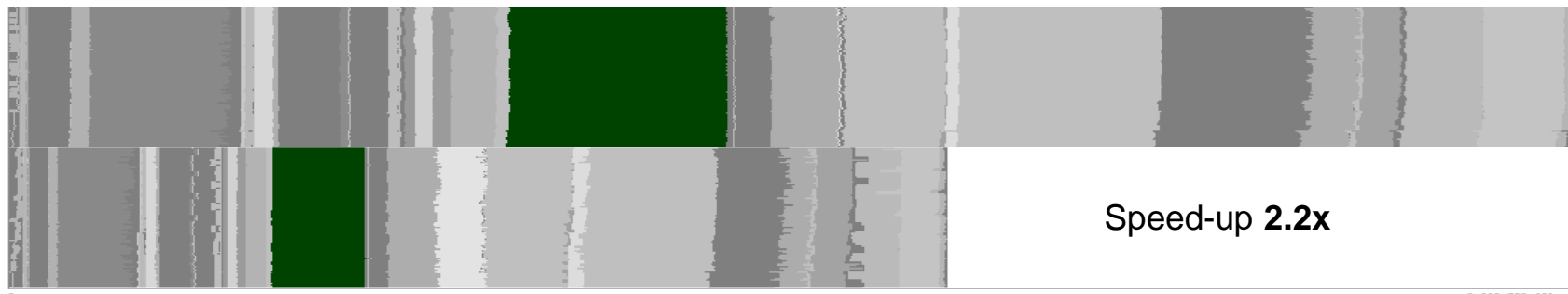
dynspg

traadv

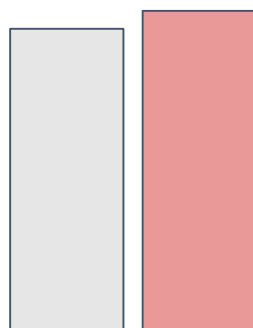
traldf

zdfgls

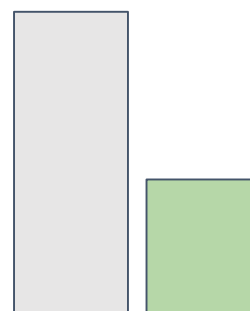
traqsr



Total
Instructions



Iteration time



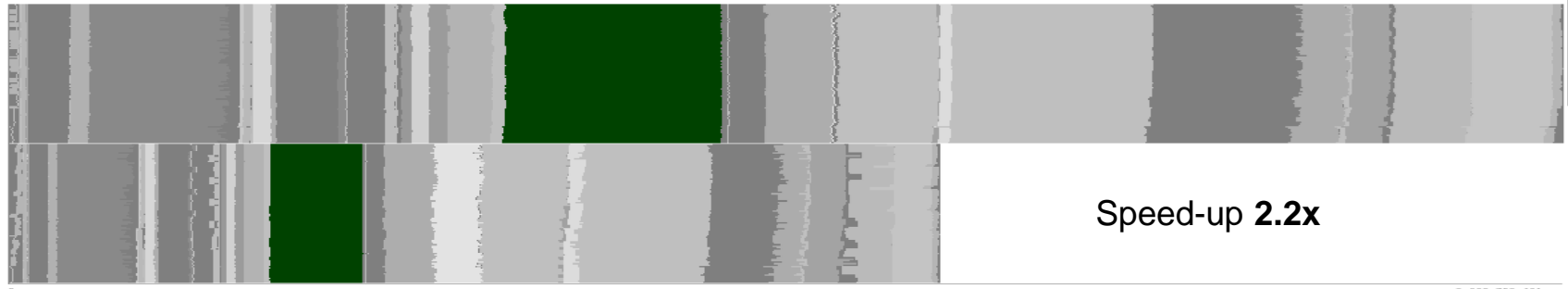
dynspg

traadv

traldf

zdfgls

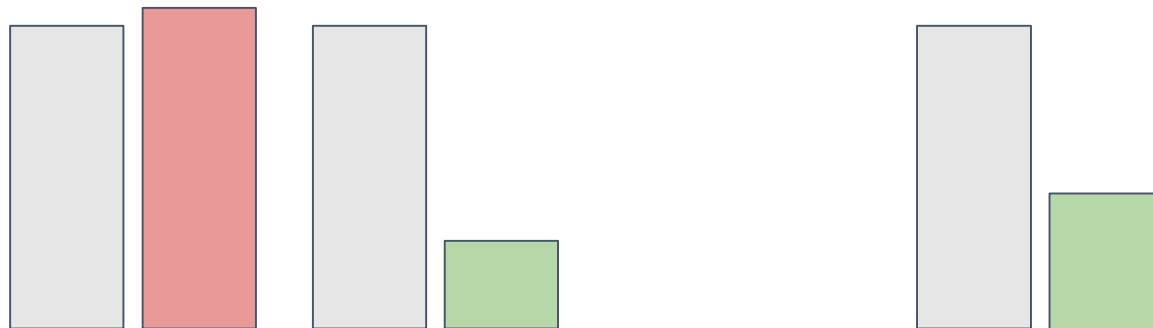
traqsr



Total
Instructions

Cycles per
Instruction

Iteration time



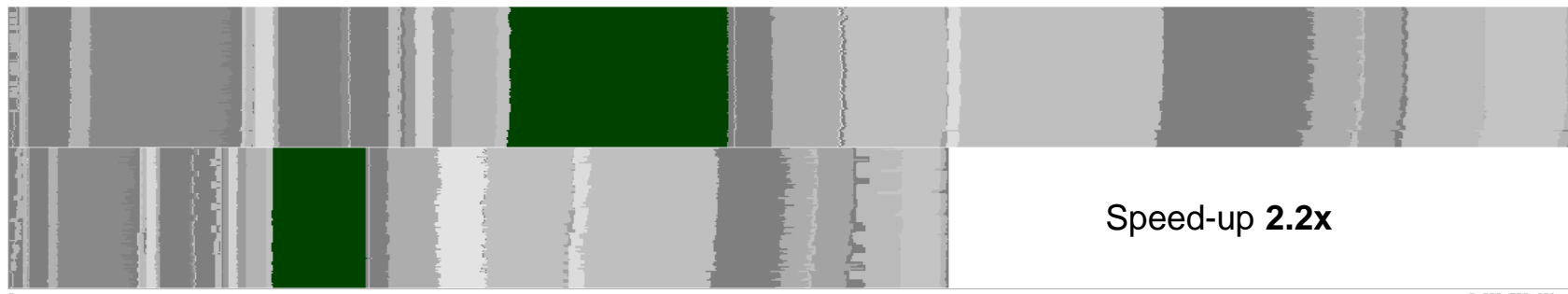
dynspg

traadv

traldf

zdfgls

traqsr

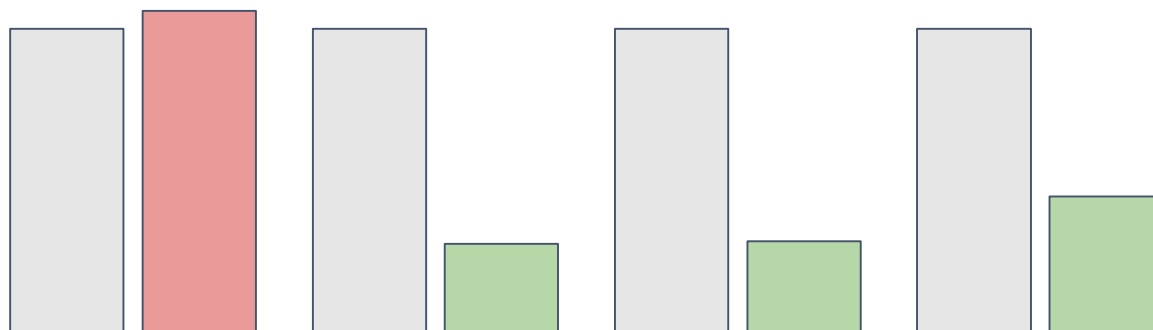


Total
Instructions

Cycles per
Instruction

L3 miss per
Instruction

Iteration time



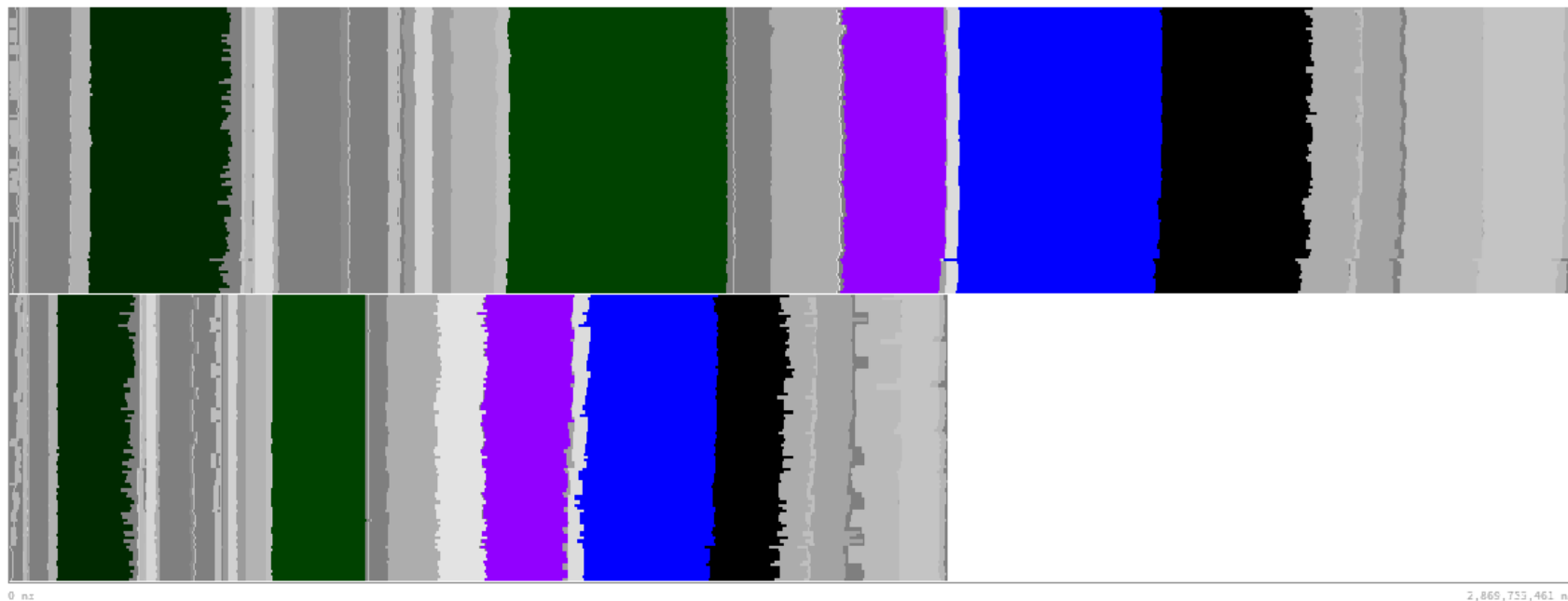
dynspg

traadv

traldf

zdfgls

traqsr



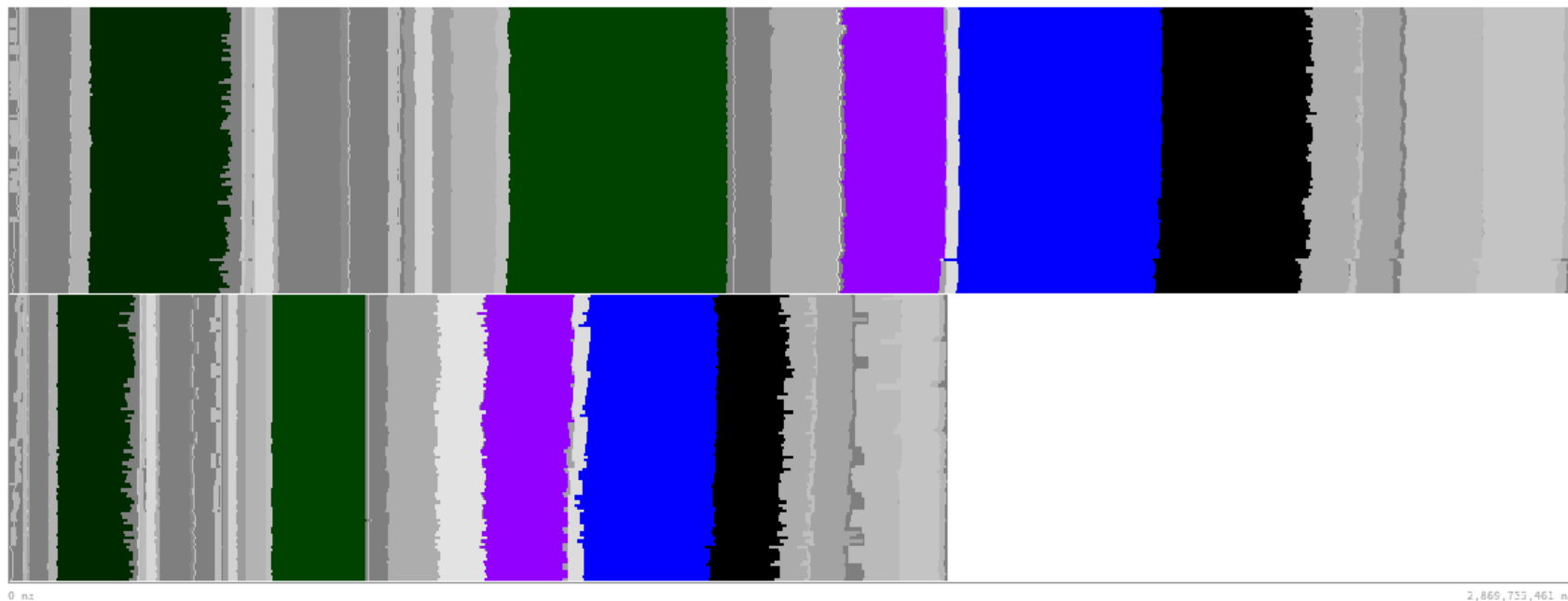
dynspg

traadv

traldf

zdfgls

traqsr



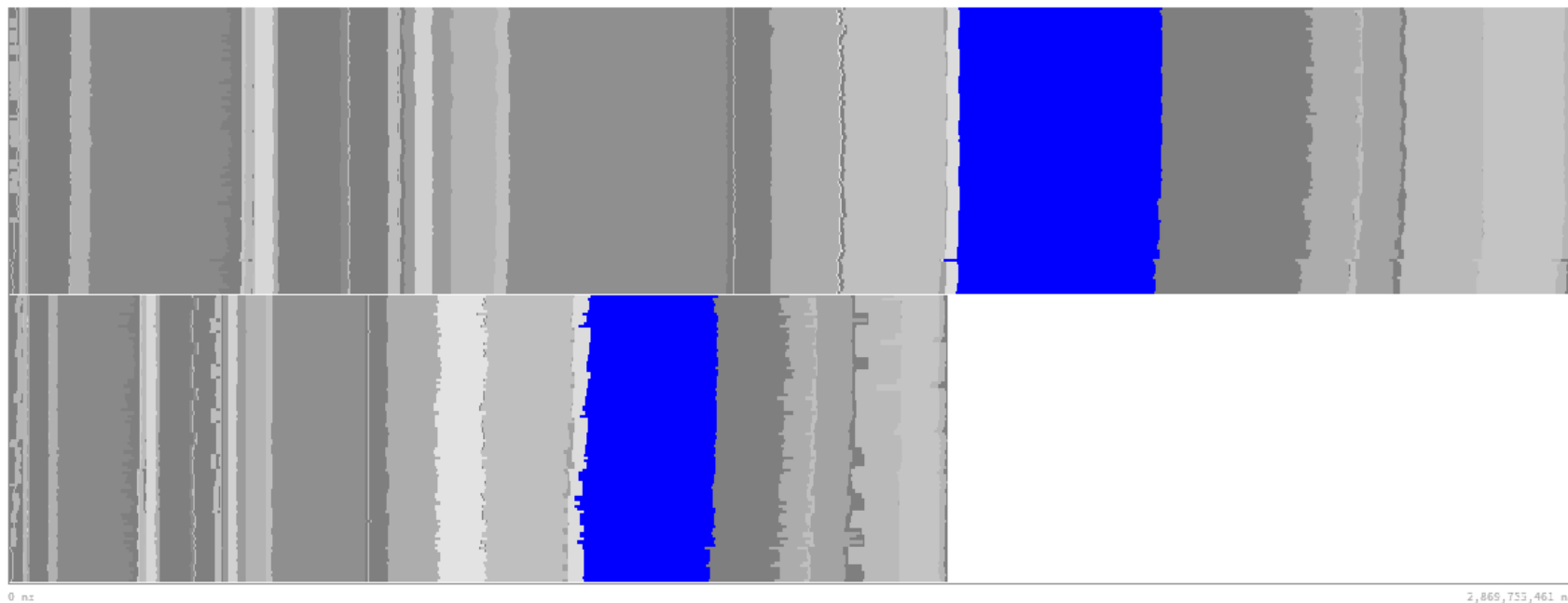
dynspg

traadv

traldf

zdfgls

traqsr



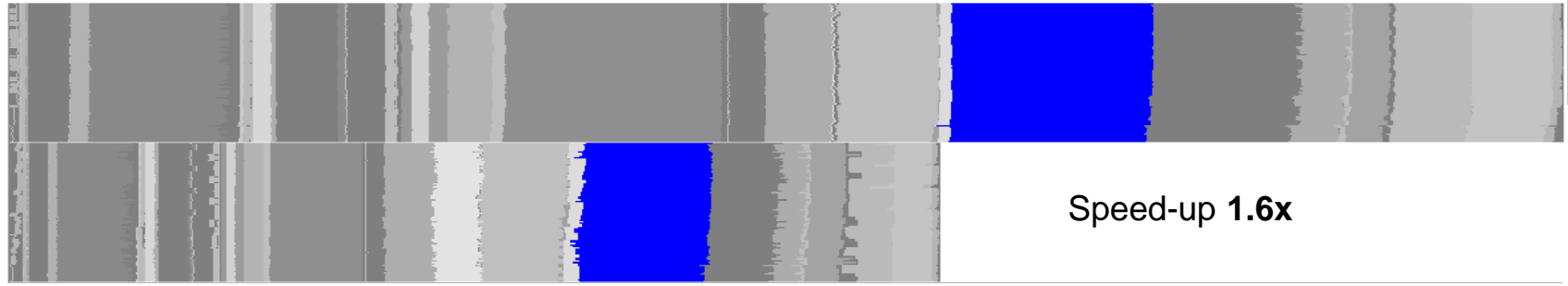
dynspg

traadv

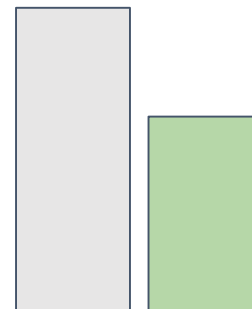
traldf

zdfgls

traqsr



Iteration time



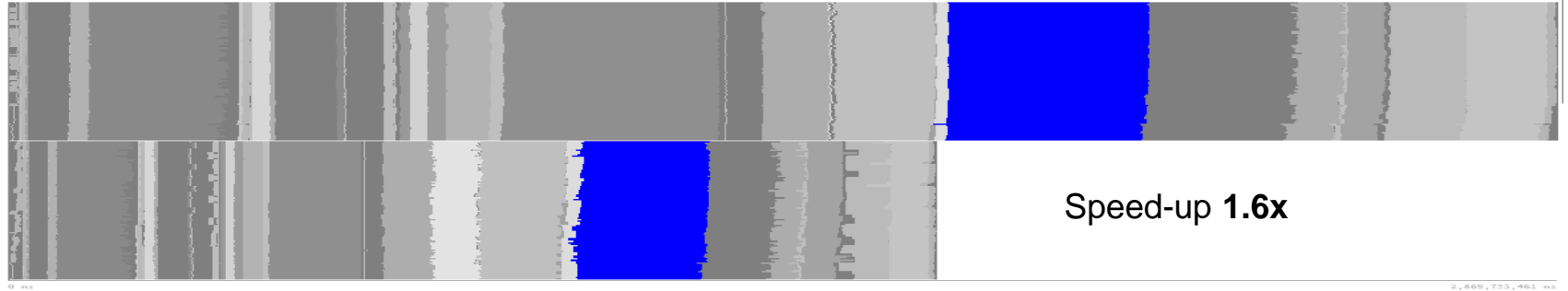
dynspg

traadv

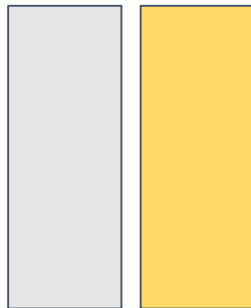
traldf

zdfgls

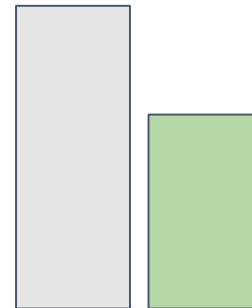
traqsr



Total
Instructions



Iteration time



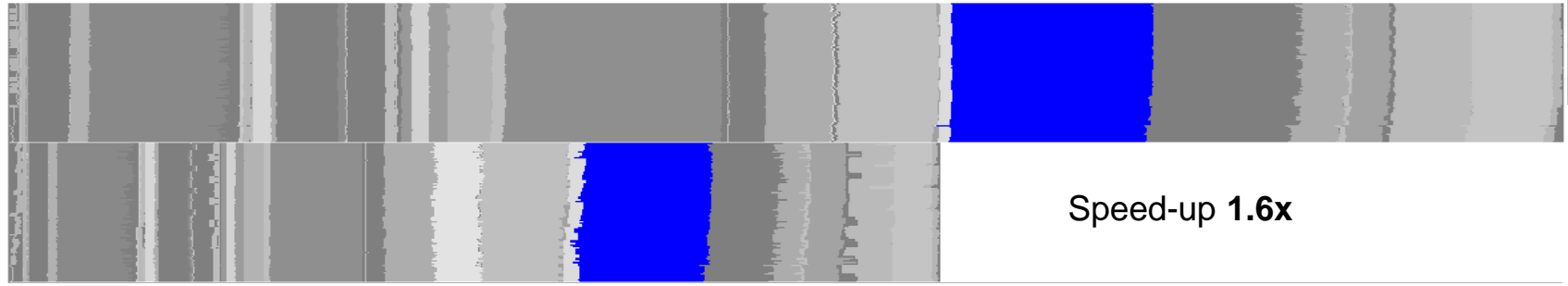
dynspg

traadv

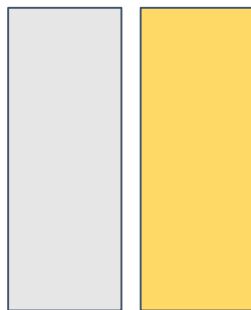
traldf

zdfgls

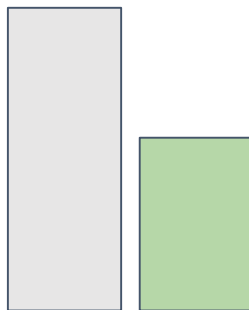
traqsr



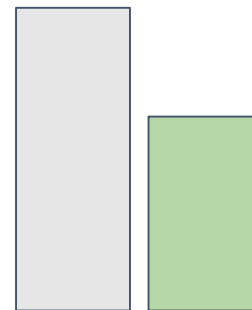
Total
Instructions



Cycles per
Instruction



Iteration time



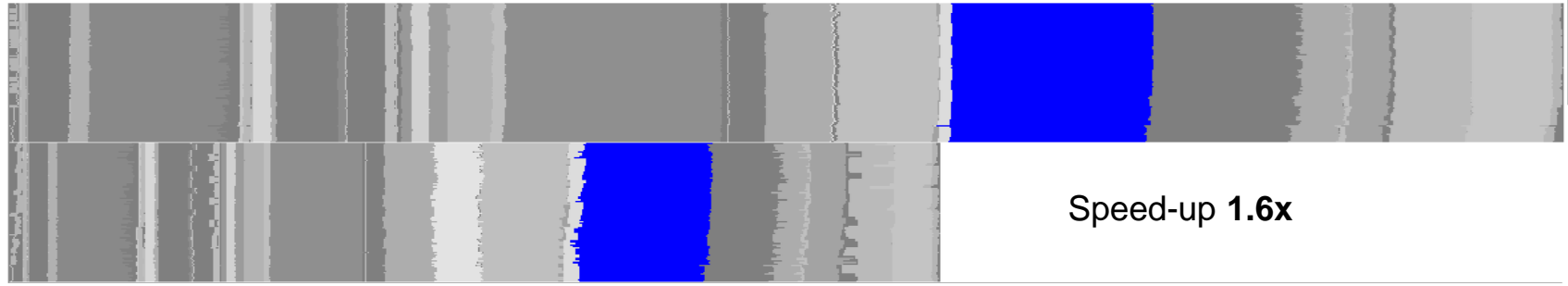
dynspg

traadv

traldf

zdfgls

traqsr

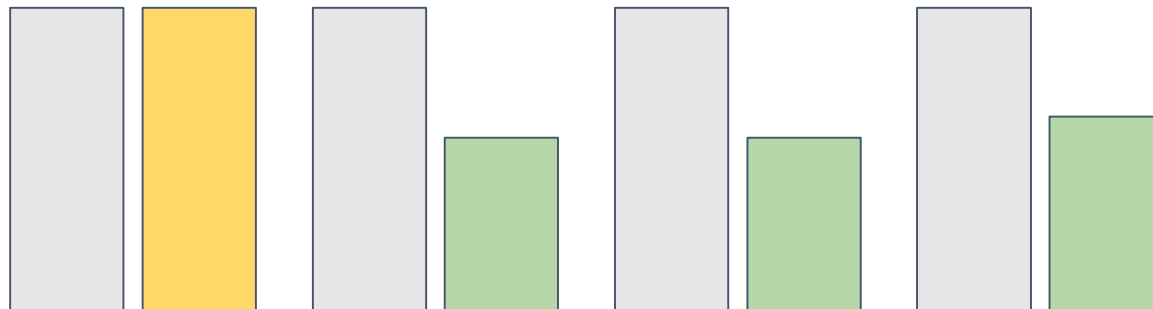


Total
Instructions

Cycles per
Instruction

L3 miss per
Instruction

Iteration time



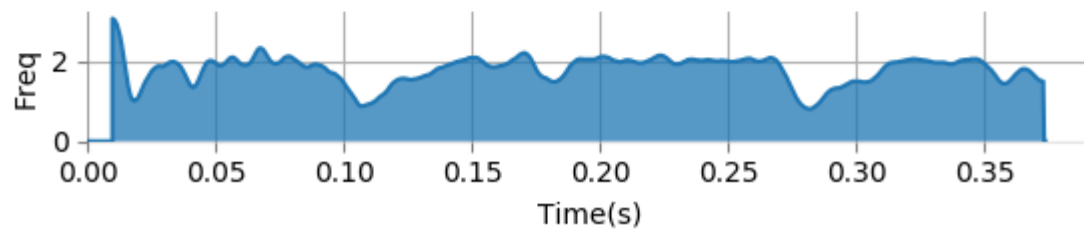
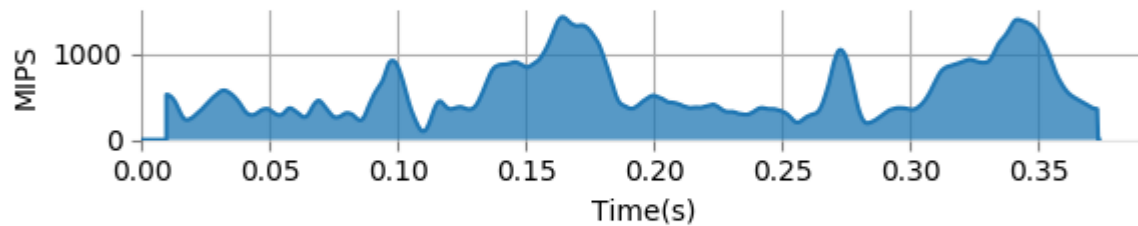
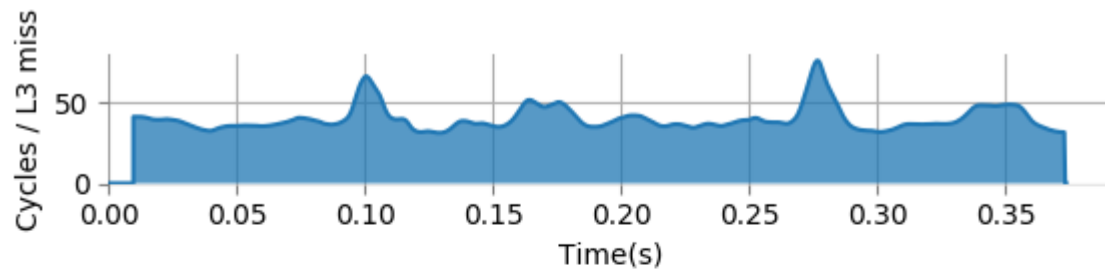
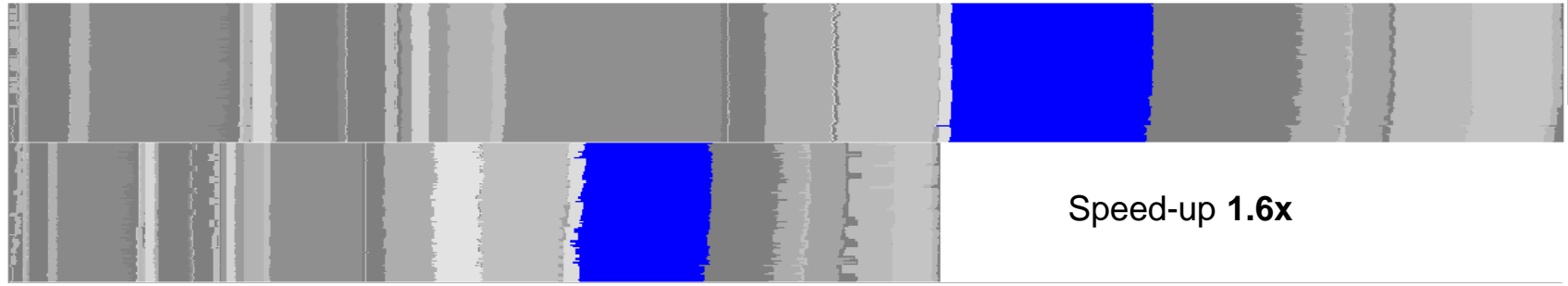
dynspg

traadv

traldf

zdfgls

traqsr



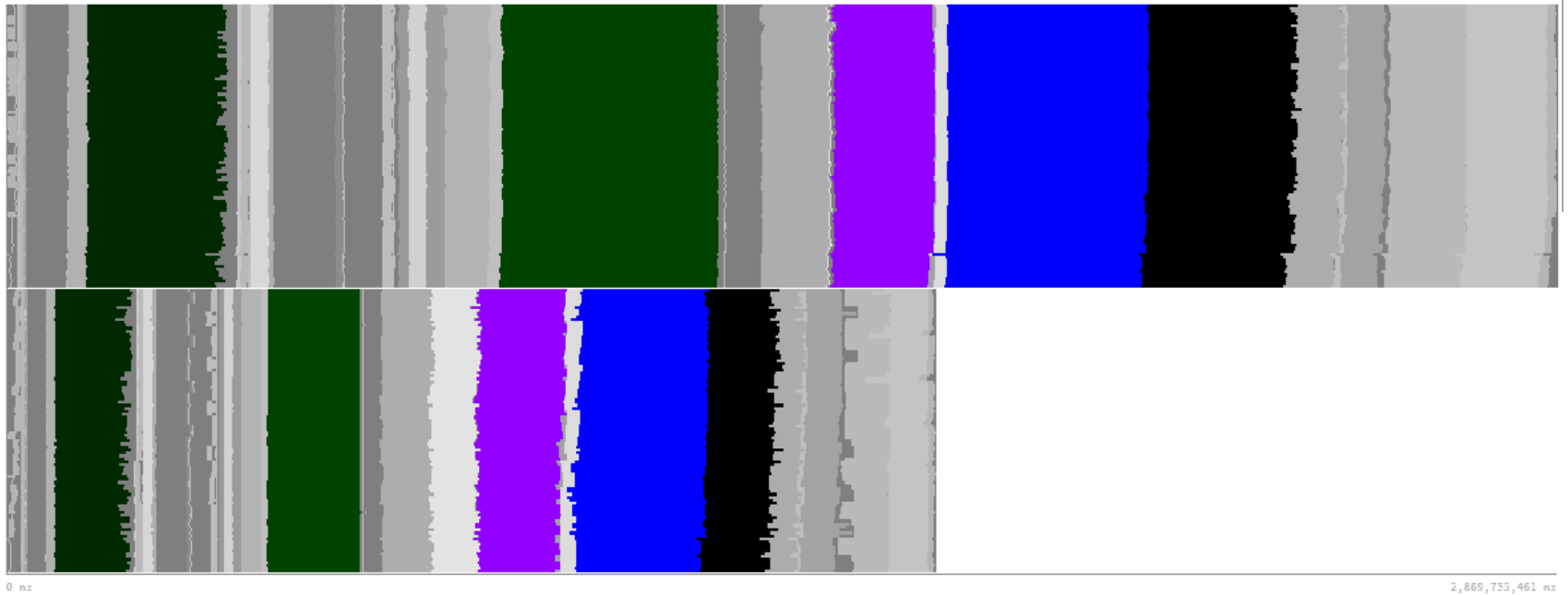
dynspg

traadv

traldf

zdfgls

traqsr



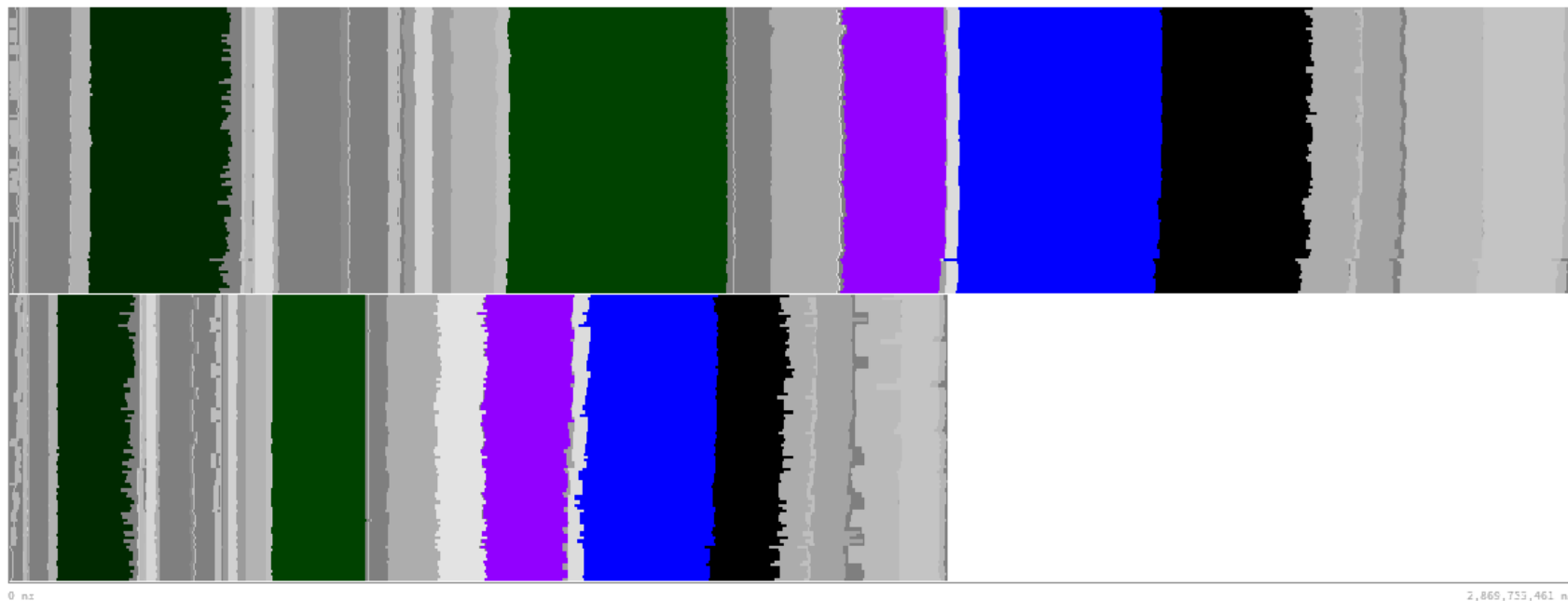
dynspg

traadv

traldf

zdfgls

traqsr



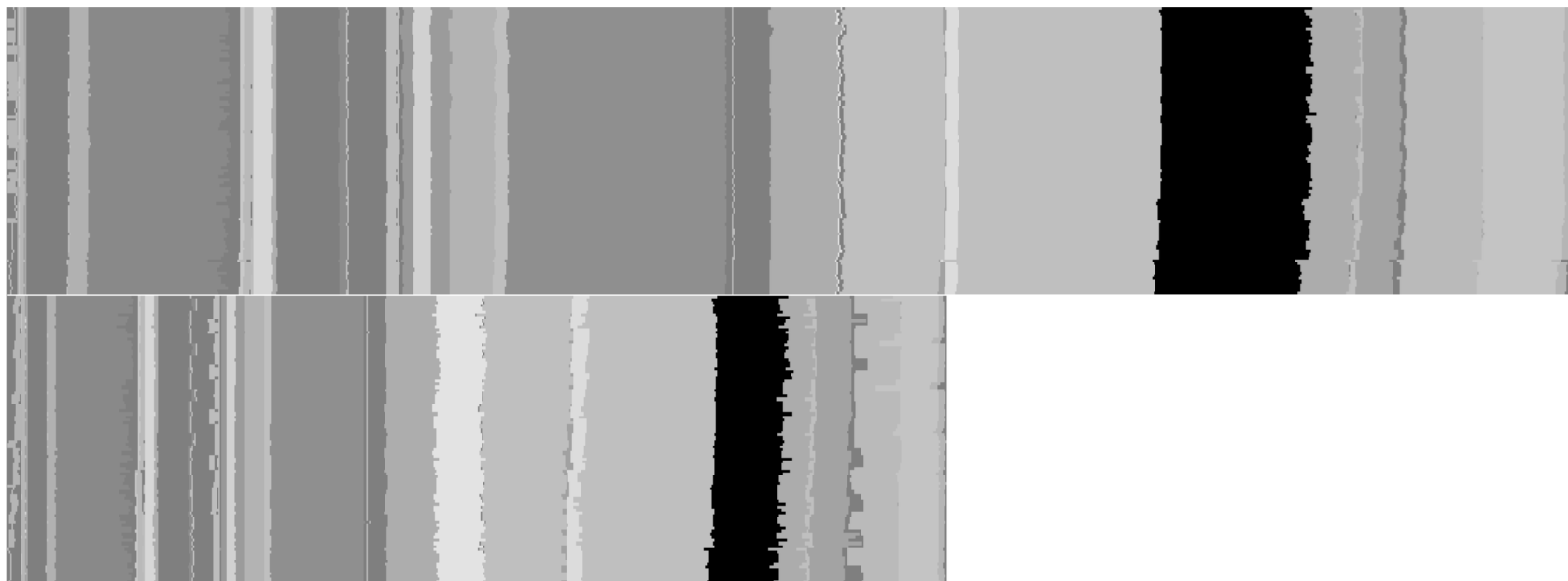
dynspg

traadv

traldf

zdfgls

traqsr



0 m

2,866,753,461 m

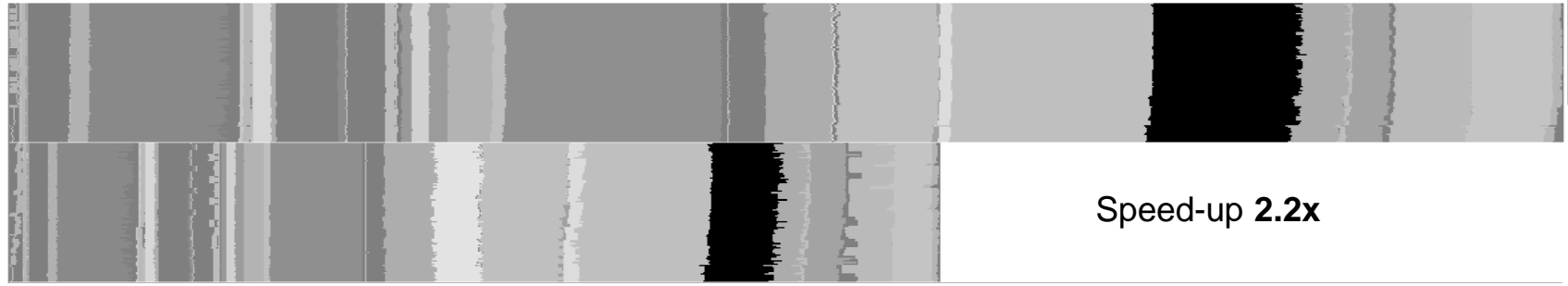
dynspg

traadv

traldf

zdfgls

traqsr

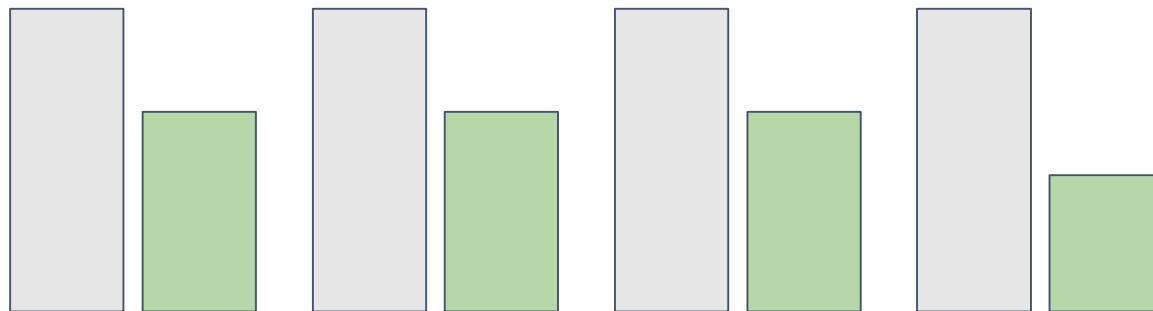


Total
Instructions

Cycles per
Instruction

L3 miss per
Instruction

Iteration time



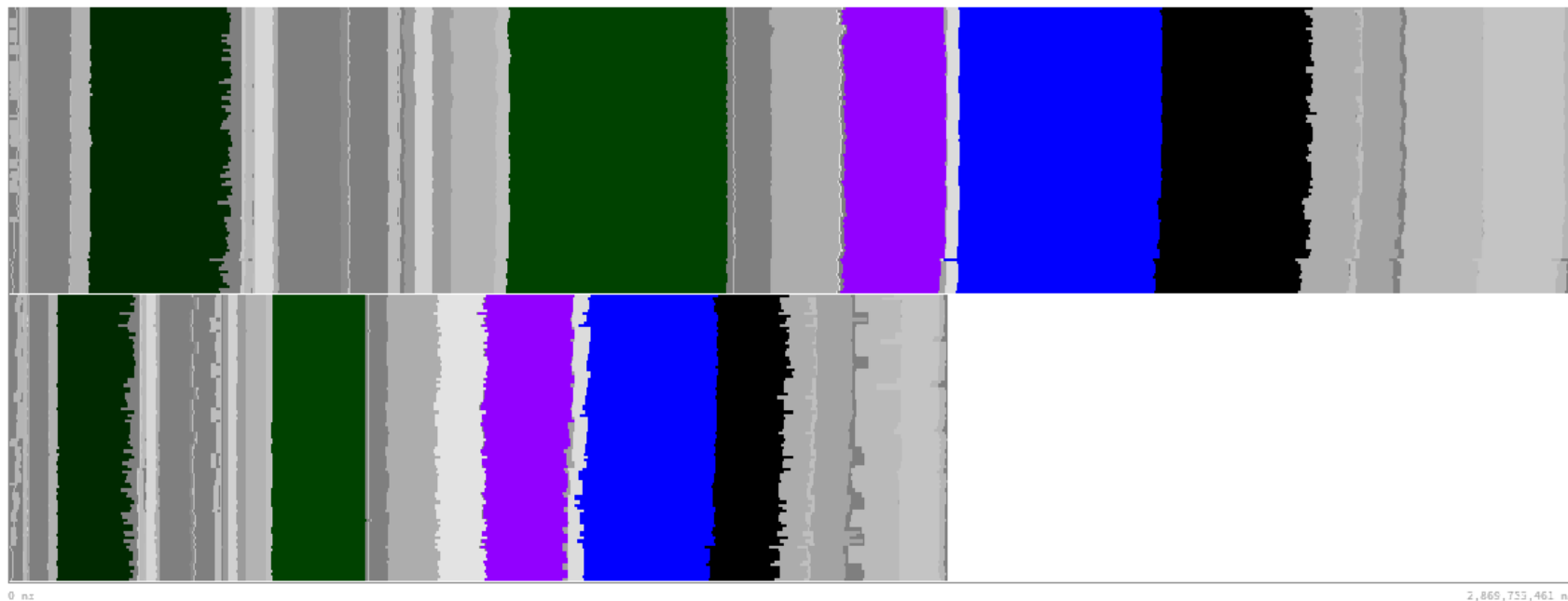
dynspg

traadv

traldf

zdfgls

traqsr



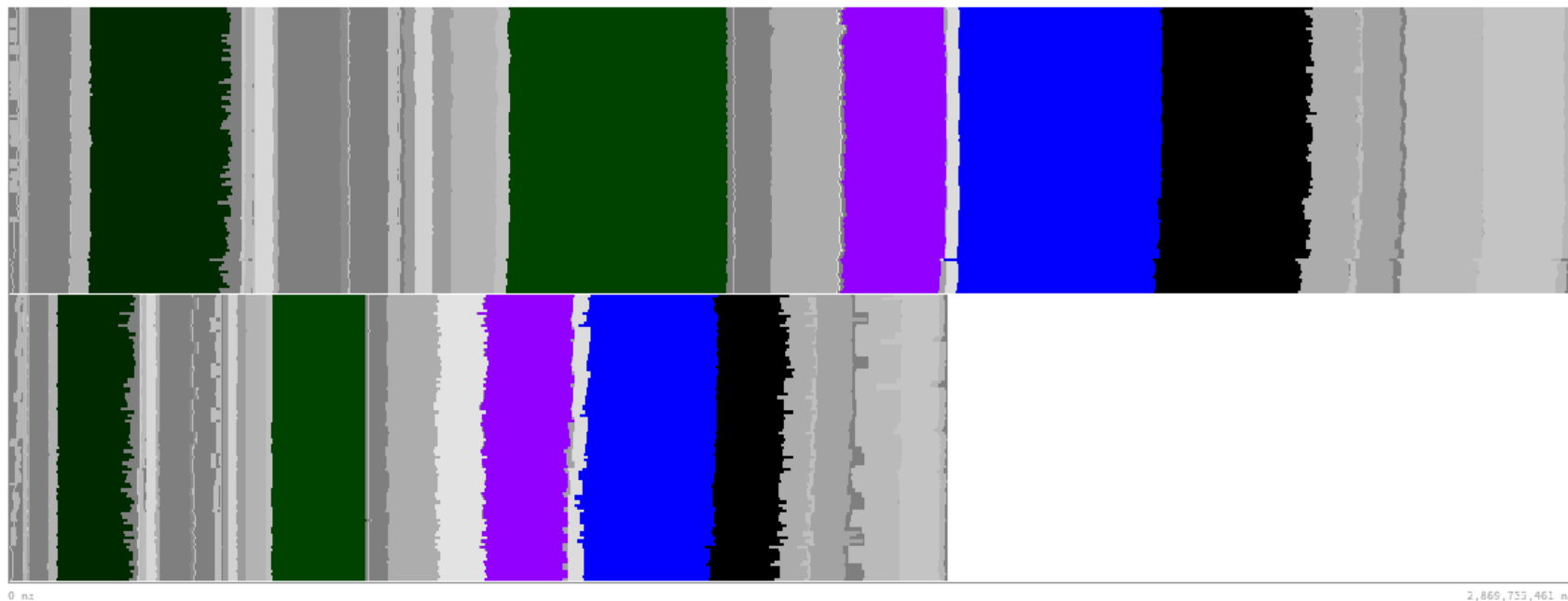
dynspg

traadv

traldf

zdfgls

traqsr



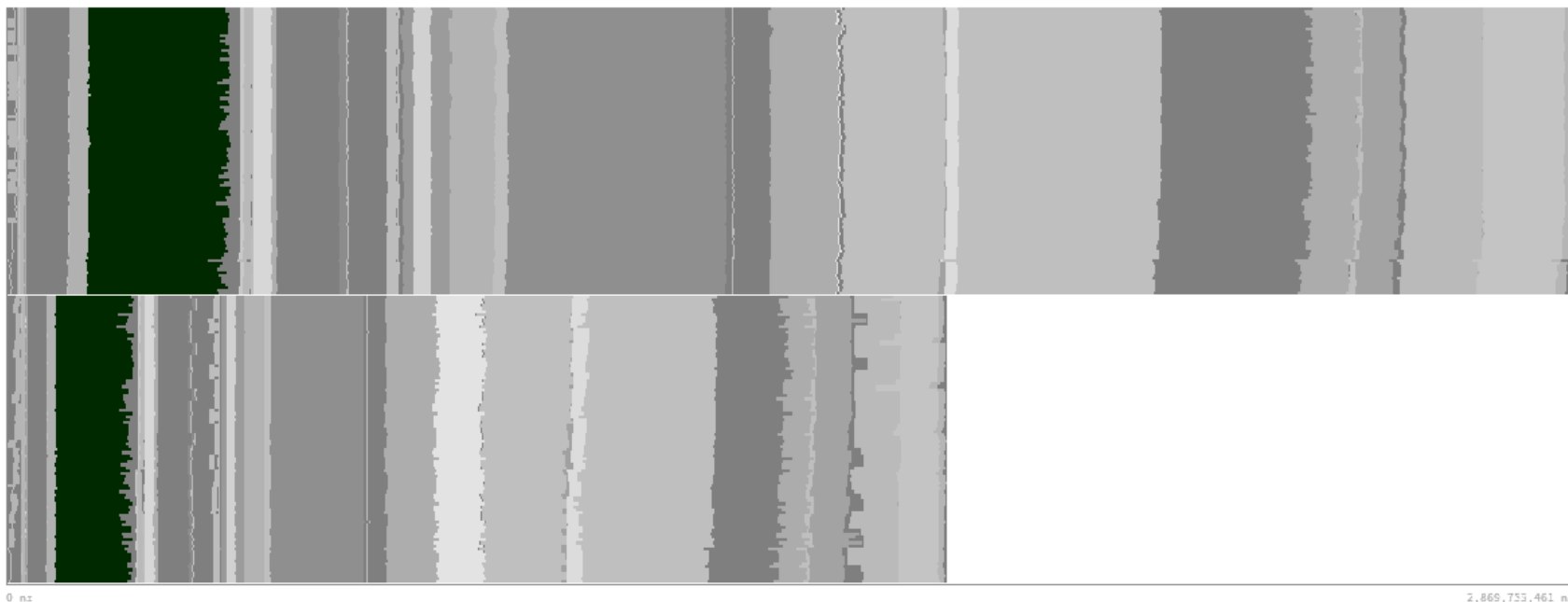
dynspg

traadv

traldf

zdfgls

traqsr



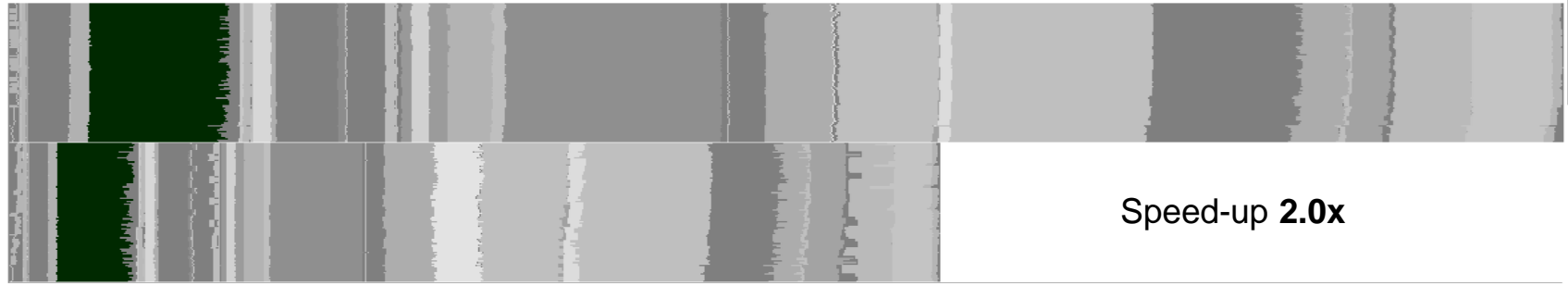
dynspg

traadv

traldf

zdfgls

traqsr

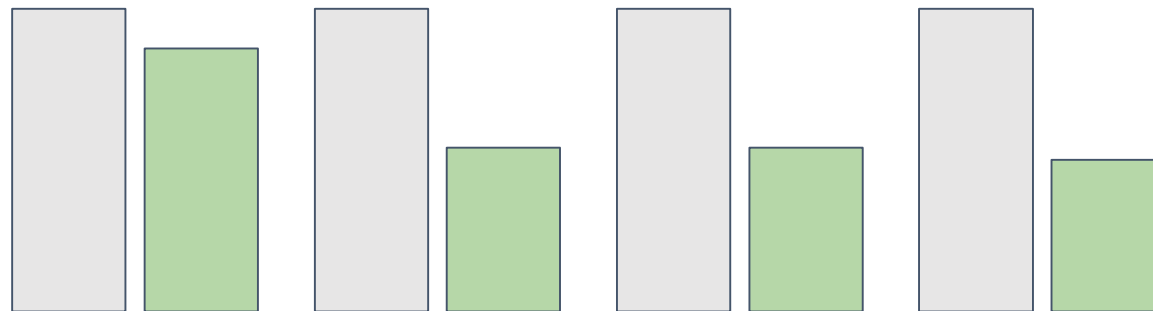


Total
Instructions

Cycles per
Instruction

L3 miss per
Instruction

Iteration time



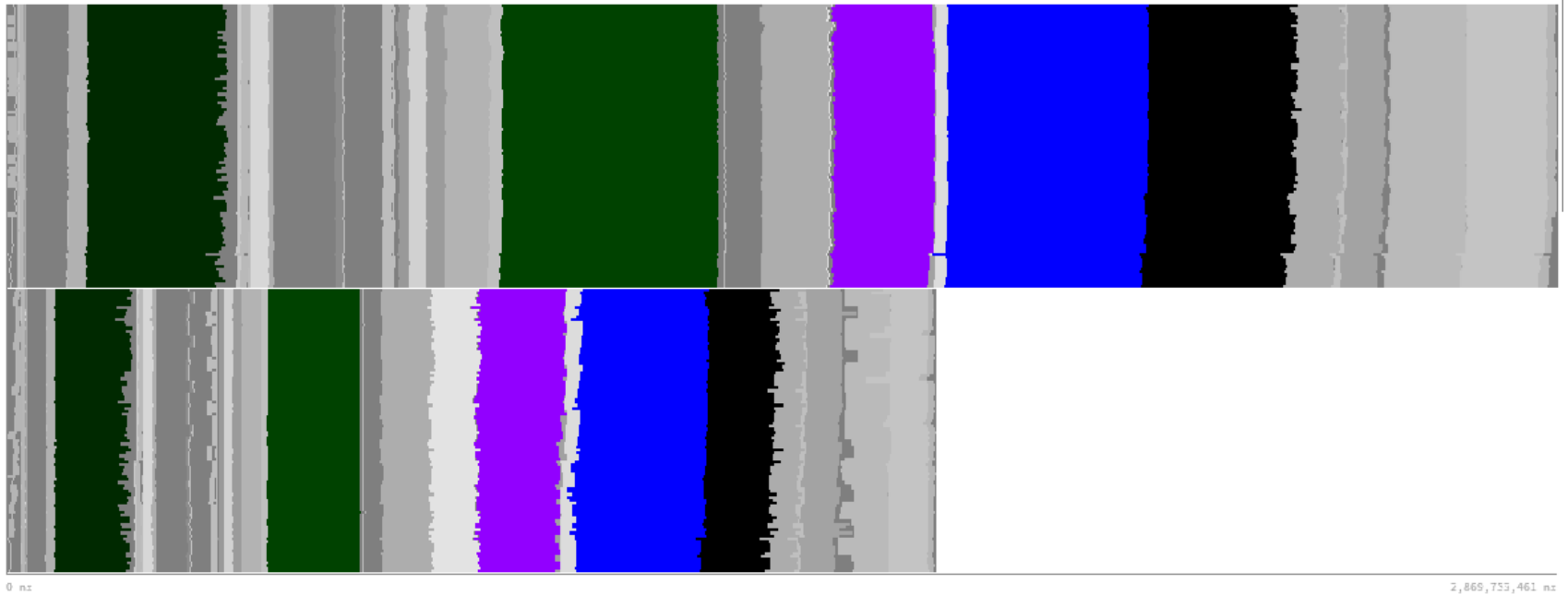
dynspg

traadv

traldf

zdfgls

traqsr



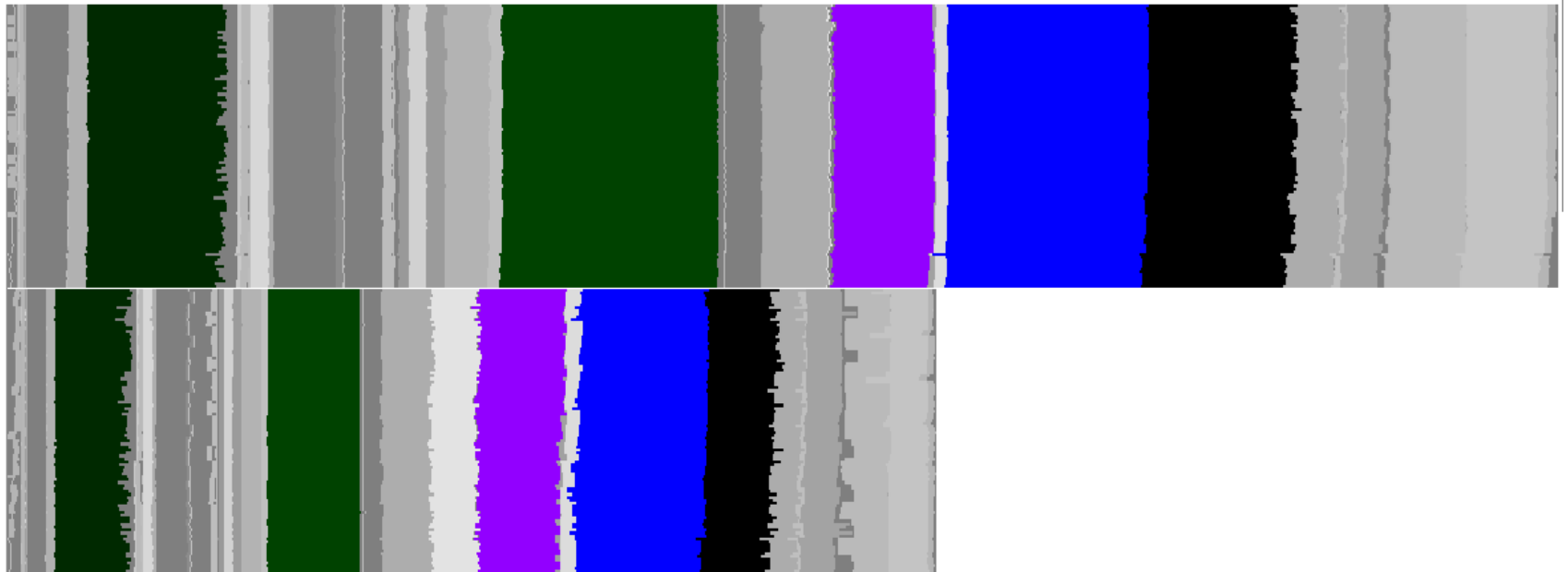
dynspg

traadv

traldf

zdfgls

traqsr



0 ns

2,868,753,461 ns

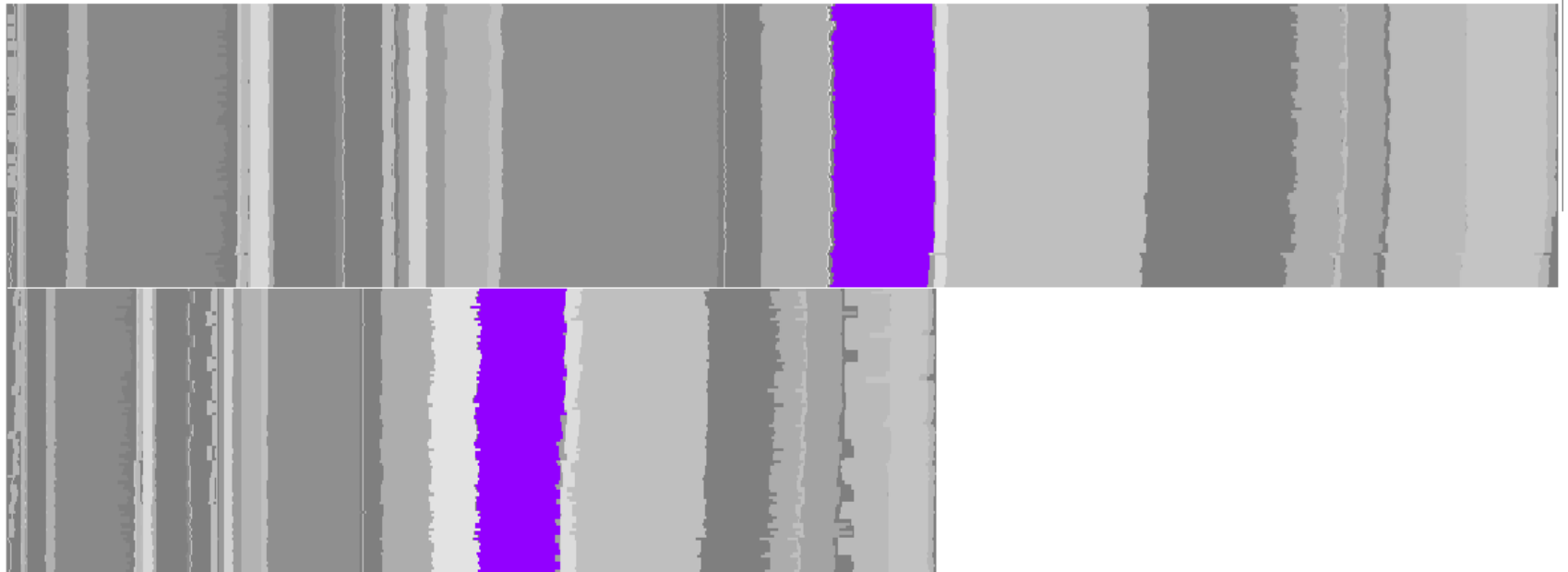
dynspg

traadv

traldf

zdfgls

traqsr



0 m

2,865,753,461 m

dynspg

traadv

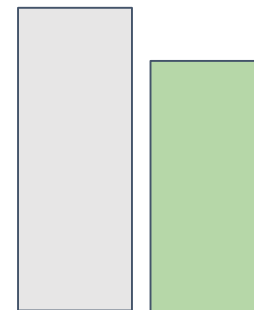
traldf

zdfgls

traqsr



Iteration time



dynspg

traadv

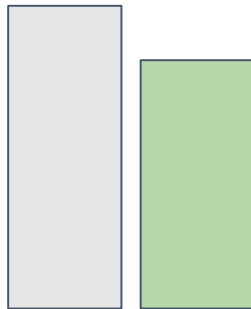
traldf

zdfgls

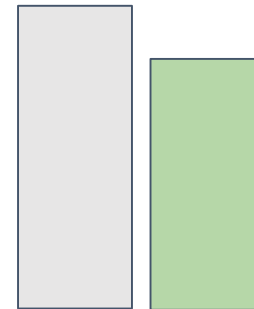
traqsr



Total
Instructions



Iteration time



dynspg

traadv

traldf

zdfgls

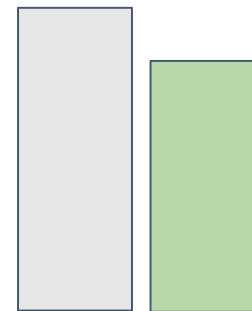
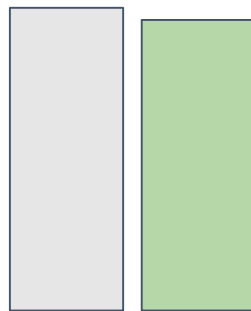
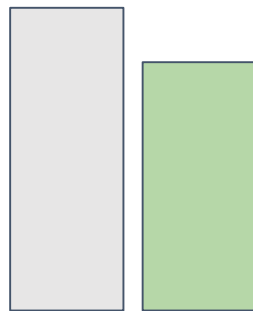
traqsr



Total
Instructions

Cycles per
Instruction

Iteration time



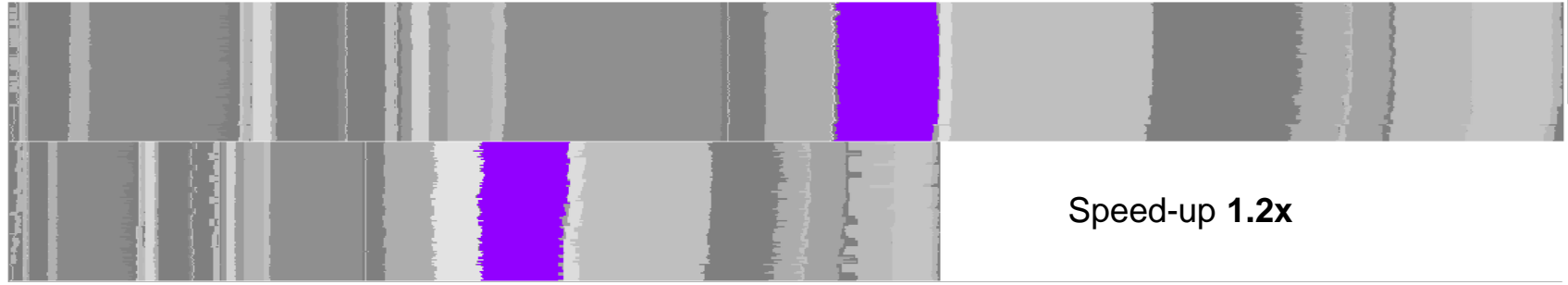
dynspg

traadv

traldf

zdfgls

traqsr

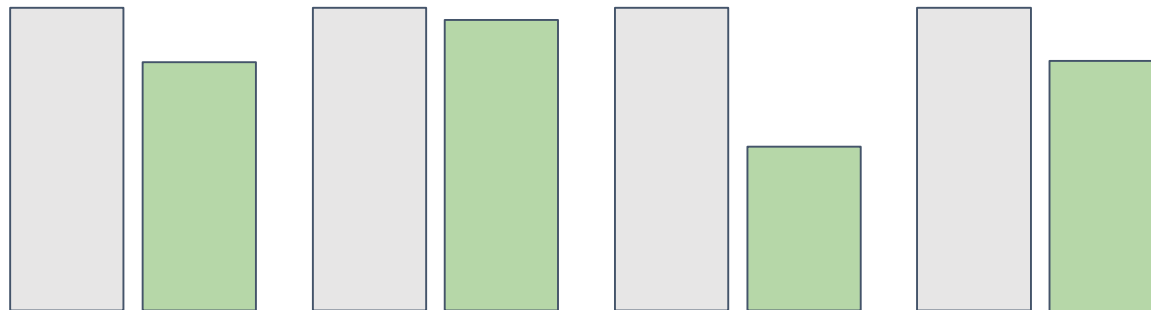


Total
Instructions

Cycles per
Instruction

L3 miss per
Instruction

Iteration time



dynspg

traadv

traldf

zdfgls

traqsr

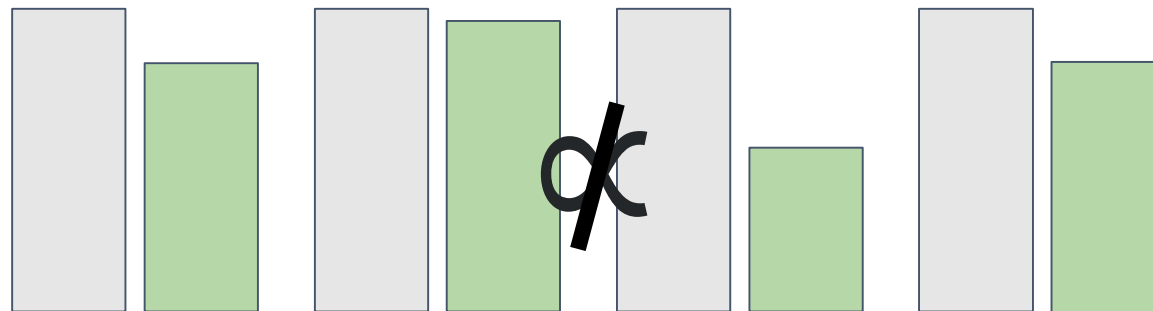


Total
Instructions

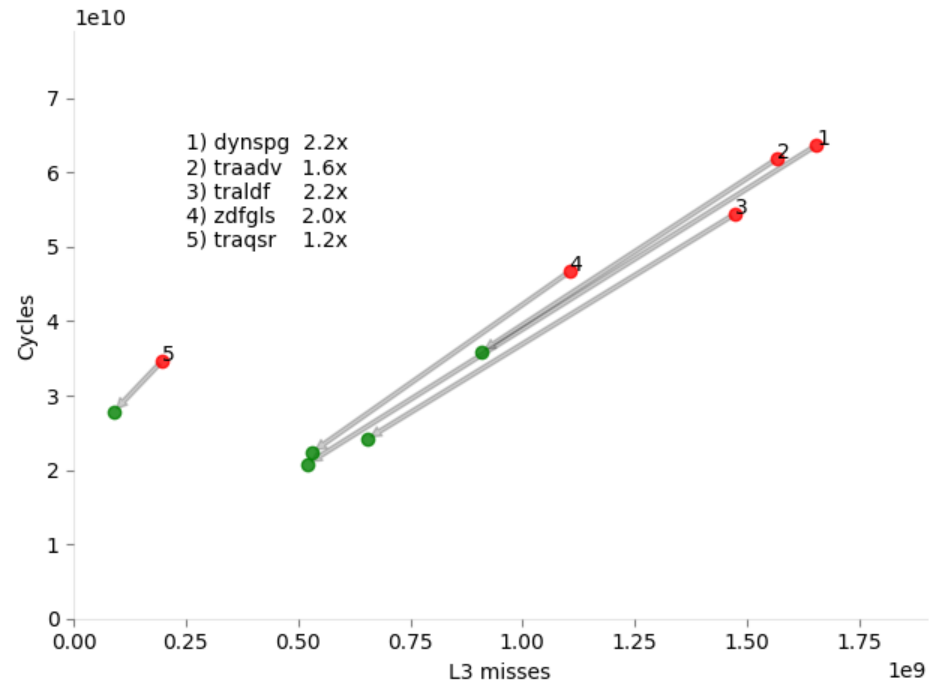
Cycles per
Instruction

L3 miss per
Instruction

Iteration time

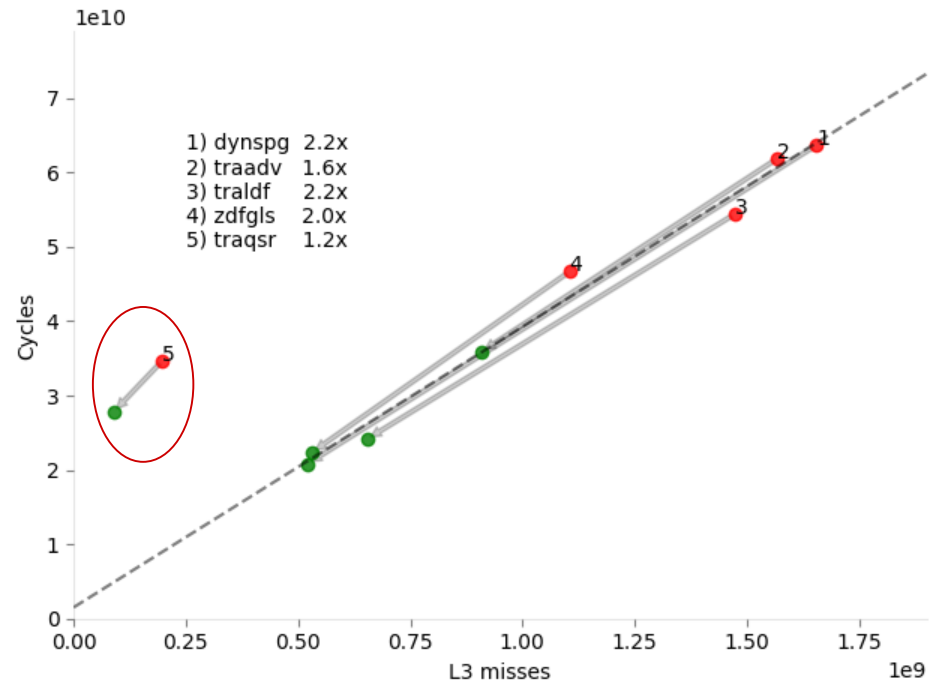


What happens with traqsr?

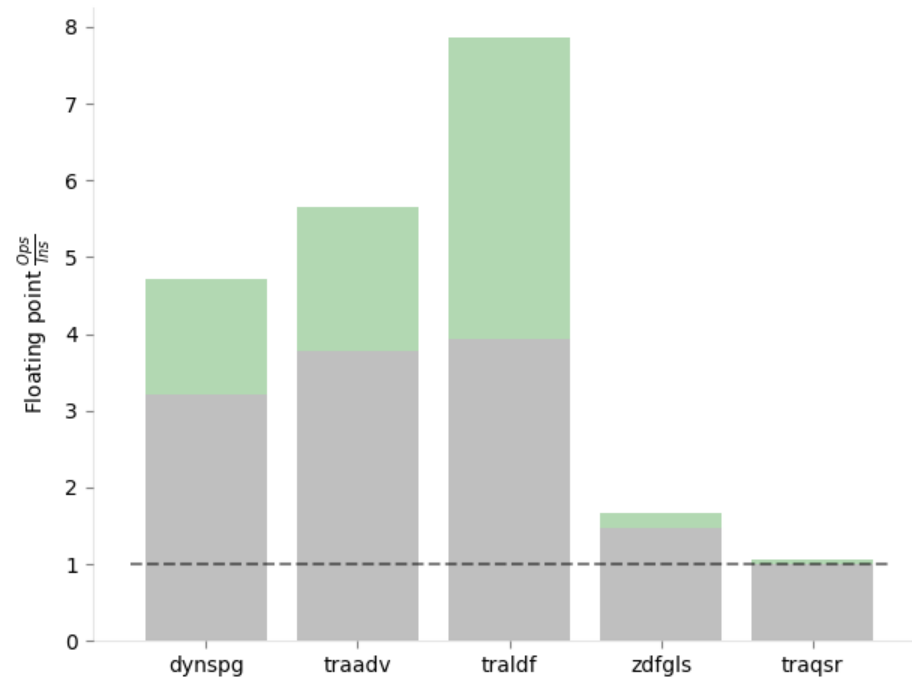


What happens with traqsr?

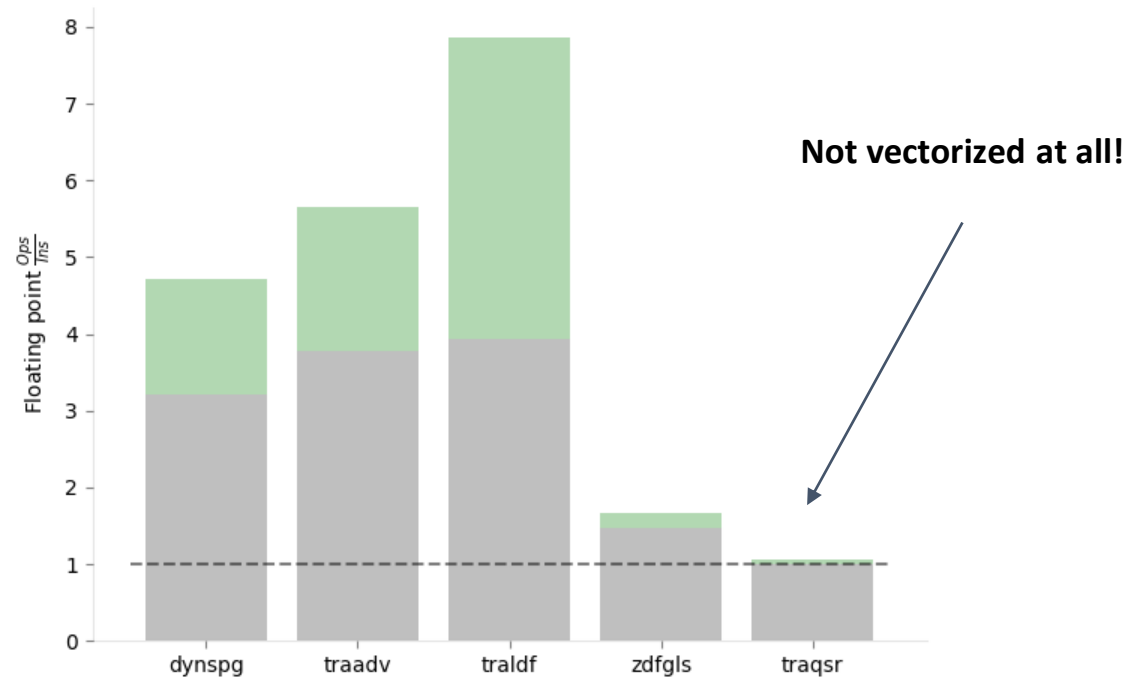
$r^2=0.99$



What happens with traqsr?

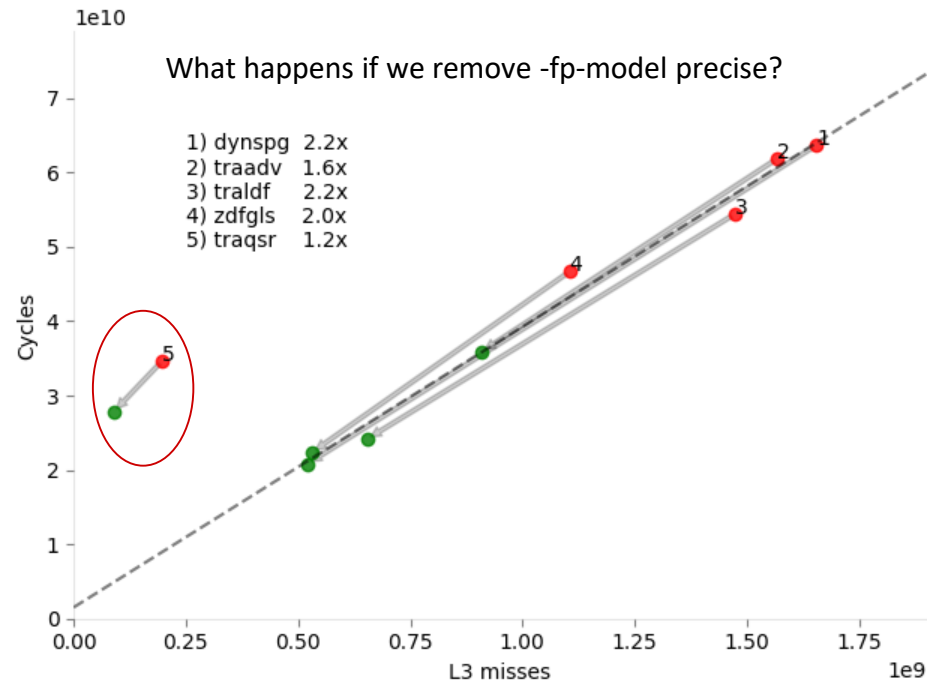


What happens with traqsr?



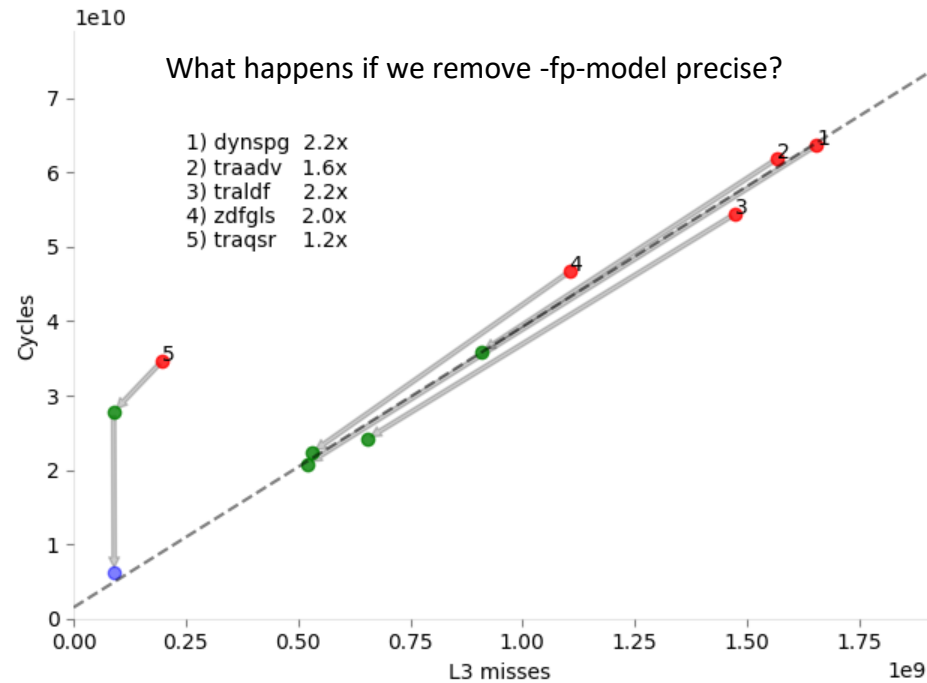
What happens with traqsr?

$r^2=0.99$

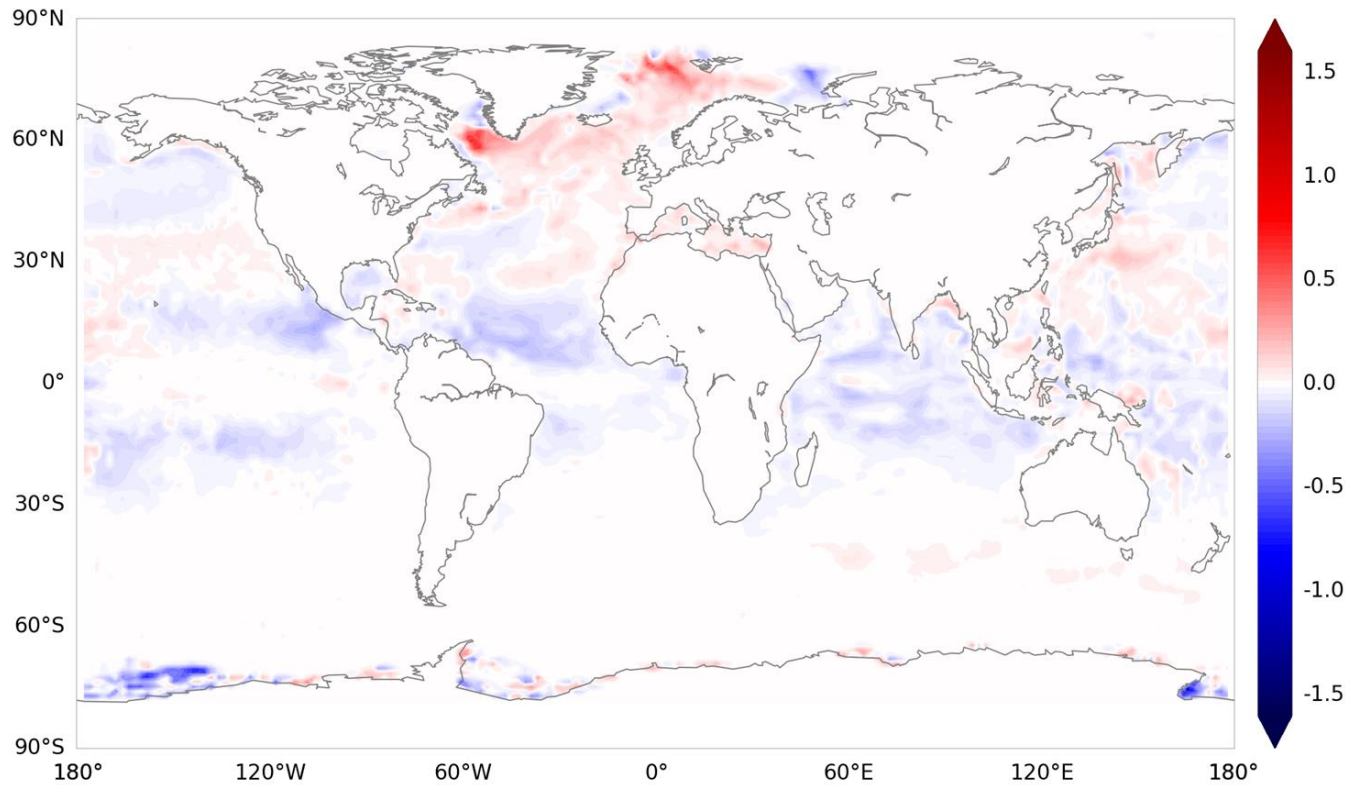


What happens with traqsr?

$r^2=0.99$

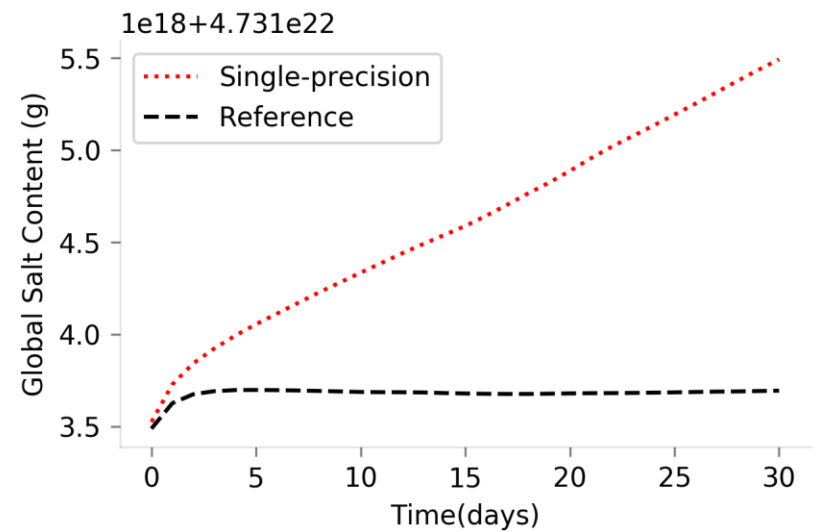
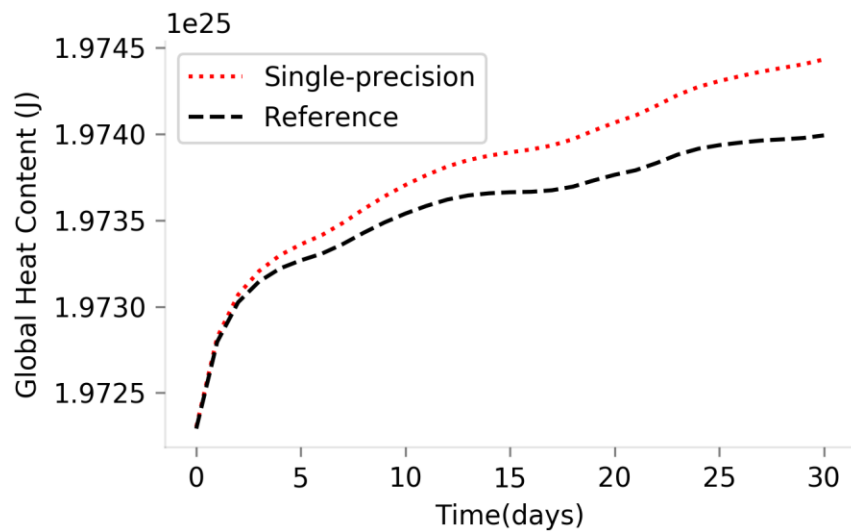


But what about the results?



Difference between double- and single-precision monthly mean of sea-surface temperature for the first month of simulation.

But what about the results?



Risks of low precision



Solution?
Precision Analysis!

source: <http://www-users.math.umn.edu/~arnold/disasters/ariane.html>

NEMO

Precision Analysis



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

Precision Analysis

- How can we find which variables need to be kept in double-precision to maintain the accuracy of the results?

Analysis Algorithm

[0 1 2 3 4 5 6 7 8 9]



Analysis Algorithm

[0 1 2 3 4 5 6 7 8 9]



[0 1 2 3 4]



[5 6 7 8 9]



Analysis Algorithm

[0 1 2 3 4 5 6 7 8 9]



[0 1 2 3 4]



[5 6 7 8 9]



[0 1 2] [3 4]



Analysis Algorithm

[0 1 2 3 4 5 6 7 8 9]



[0 1 2 3 4] 

[5 6 7 8 9] 

[0 1 2]  [3 4] 

[3]⓪ [4] ⓪

Analysis Algorithm

[0 1 2 3 4 5 6 7 8 9]



[0 1 2 3 4] 

[5 6 7 8 9] 

[0 1 2]  [3 4] 

[3]  [4] 

Analysis Algorithm

[0 1 2 3 4 5 6 7 8 9]



[0 1 2 3 4]

[5 6 7 8 9]

[0 1 2] [3 4]

[3] [4]

Analysis Algorithm

[0 1 2 3 4 5 6 7 8 9]



[0 1 2 3 4] ?

[5 6 7 8 9] ✓

[0 1 2] ✓ [3 4] ✓

[3] ✓ [4] ✗

Analysis Algorithm

[0 1 2 3 4 5 6 7 8 9]



[0 1 2 3 4] ✓

[5 6 7 8 9] ✓

[0 1 2] ✓ [3 4] ✓

[3] ✓ [4] ✗

Analysis Algorithm

[0 1 2 3 4 5 6 7 8 9]



[0 1 2 3 4] ✓

[5 6 7 8 9] ✓

[0 1 2] ✓ [3 4] ✓

[3] ✓ [4] ✗

Analysis Algorithm

[0 1 2 3 4 5 6 7 8 9]



[0 1 2 3 4] ✓

[5 6 7 8 9] ✓

[0 1 2] ✓ [3 4] ✓

[3] ✓ [4] ✗



Variable 4 must be kept in double-precision.

Analysis Algorithm

[0 1 2 3 4 5 6 7 8 9]



[0 1 2 3 4] ?

[5 6 7 8 9] ✓

[0 1 2] ✓ [3 4] ✓

[3] ✓ [4] ✗

Analysis Algorithm

[0 1 2 3 4 5 6 7 8 9]



[0 1 2 3 4]



[5 6 7 8 9]



[0 1 2]



[3 4]



[3]



[4]



Analysis Algorithm

[0 1 2 3 4 5 6 7 8 9]



[0 1 2 3 4]



[5 6 7 8 9]



[0 1 2] ✓ [?] > [3 4] ✓

[3] ✓ [4] ✗

Analysis Algorithm

[0 1 2 3 4 5 6 7 8 9]



[0 1 2 3 4] ?

[5 6 7 8 9] ✓

[0 1 2]+[3] ? [3 4] ✓

[0 1]+[3] ? [2]+[3] ? [3] ✓ [4] ✗

Analysis Algorithm

[0 1 2 3 4 5 6 7 8 9]



[0 1 2 3 4] ✓

[5 6 7 8 9] ✓

[0 1 2]+[3] ✓ [3 4] ✓

[0 1]+[3] ✓ [2]+[3] ✗ [3] ✓ [4] ✗

Analysis Algorithm

[0 1 2 3 4 5 6 7 8 9]



[0 1 2 3 4] ✓

[5 6 7 8 9] ✓

[0 1 2]+[3] ✓ [3 4] ✓

[0 1]+[3] ✓ [2]+[3] ✗ [3] ✓ [4] ✗

Analysis Algorithm

[0 1 2 3 4 5 6 7 8 9]



[0 1 2 3 4] ✓

[5 6 7 8 9] ✓

[0 1 2]+[3] ✓ [3 4] ✓

[0 1]+[3] ✓ [2]+[3] ✗ [3] ✓ [4] ✗



Variables 2 and 4 must be kept in double-precision.

More info: How to use mixed precision in ocean models. <https://www.geosci-model-dev-discuss.net/gmd-2019-20/>

Precision Analysis

- How can we find which variables need to be kept in double-precision to maintain the accuracy of the results?:
 - How can we measure the effect of reducing the precision of an arbitrary set of variables?
 - How can we verify the accuracy of the results?

Reduced Precision Emulator

Overview

The library contains a derived type: `rpe_var`. This type can be used in place of real-valued variables to perform calculations with floating-point numbers represented with a reduced number of bits in the floating-point significand.

Basic use of the reduced-precision type

The `rpe_var` type is a simple container for a double precision floating point value. Using an `rpe_var` instance is as simple as declaring it and using it just as you would a real number:

```
TYPE(rpe_var) :: myvar  
  
myvar = 12  
myvar = myvar * 1.287  ! reduced-precision result is stored in `myvar`
```

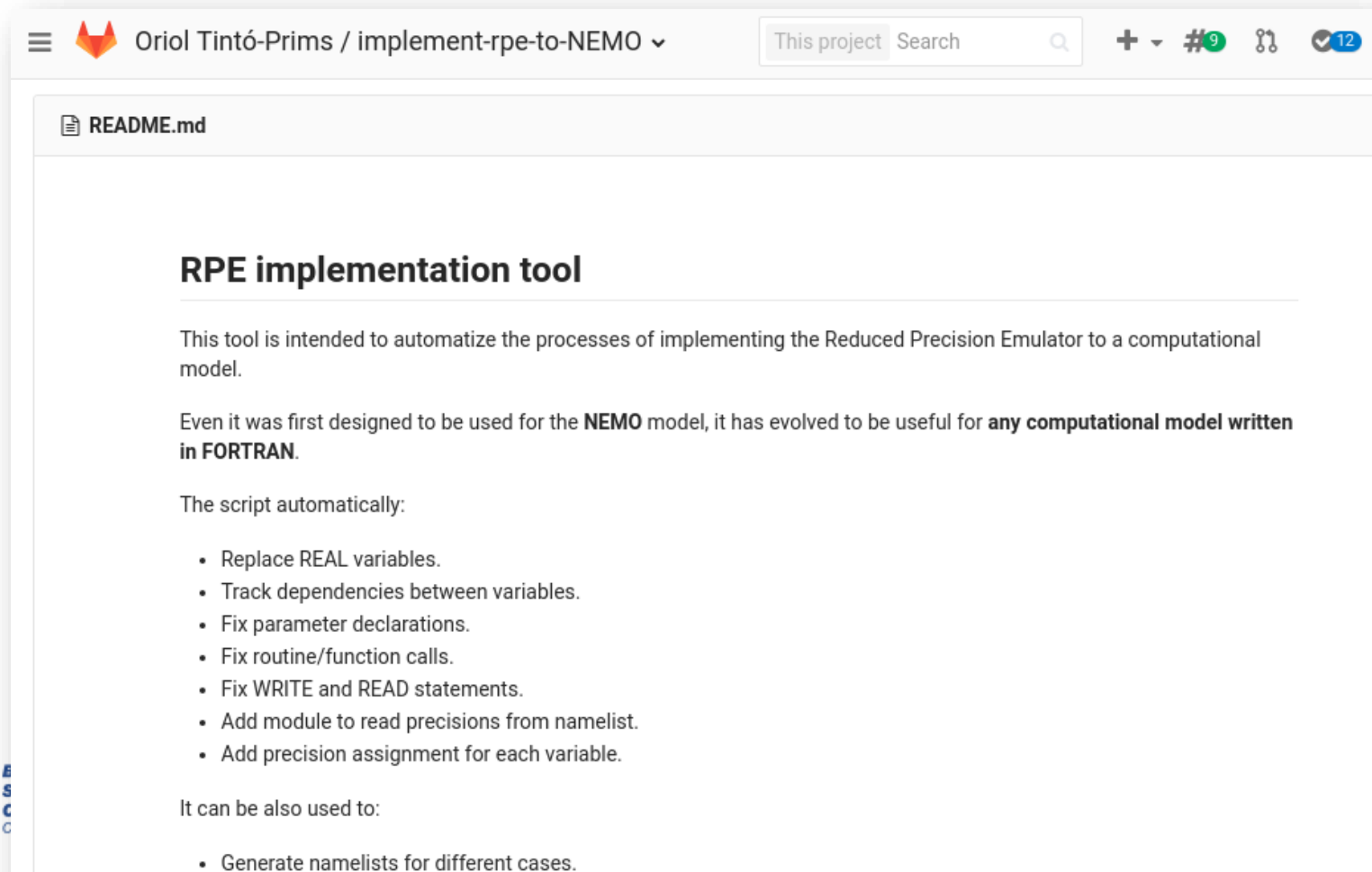
Controlling the precision

The precision used by reduced precision types can be controlled at two different levels. Each reduced precision variable has an `sbits` attribute which controls the number of explicit bits in its significand. This can be set independently for different variables, and comes into effect after it is explicitly set.

```
TYPE(rpe_var) :: myvar1  
TYPE(rpe_var) :: myvar2  
  
! Use 16 explicit bits in the significand of myvar1, but only 12 in the  
! significand of myvar2.  
myvar1%sbits = 16  
myvar2%sbits = 12
```


Implementing the emulator

- A Python tool to automate the implementation process was created.



The screenshot shows a GitHub repository page for 'Oriol Tintó-Prims / implement-rpe-to-NEMO'. The repository name is visible in the header. Below the header, the 'README.md' file is open, displaying the following content:

RPE implementation tool

This tool is intended to automatize the processes of implementing the Reduced Precision Emulator to a computational model.

Even it was first designed to be used for the **NEMO** model, it has evolved to be useful for **any computational model written in FORTRAN**.

The script automatically:

- Replace REAL variables.
- Track dependencies between variables.
- Fix parameter declarations.
- Fix routine/function calls.
- Fix WRITE and READ statements.
- Add module to read precisions from namelist.
- Add precision assignment for each variable.

It can be also used to:

- Generate namelists for different cases.

In the bottom left corner, there is a BSC logo. In the bottom right corner, there is a 'vace' logo with the text 'IN SIMULATION OF WEATHER' below it.

RPE in NEMO:

What we can do with it?

With a **single binary**, we can specify the number of significant bits used for each real variable declaration within the code through a **namelist**.

```

1  ! namelist variable precisions
2  $precisions
3  emulator_variable_precisions(1) = 10  ! Variable:      ad_u Routine:      ad_balance_tile Module:      ad_balance
4  emulator_variable_precisions(2) = 10  ! Variable:      ad_v Routine:      ad_balance_tile Module:      ad_balance
5  emulator_variable_precisions(3) = 10  ! Variable:      ad_zeta Routine:      ad_balance_tile Module:      ad_balance
6  emulator_variable_precisions(5) = 10  ! Variable:      pc_r2d Routine:      ad_balance_tile Module:      ad_balance
7  emulator_variable_precisions(6) = 10  ! Variable:      r_r2d Routine:      ad_balance_tile Module:      ad_balance
8  emulator_variable_precisions(7) = 10  ! Variable:      br_r2d Routine:      ad_balance_tile Module:      ad_balance
9  emulator_variable_precisions(8) = 10  ! Variable:      p_r2d Routine:      ad_balance_tile Module:      ad_balance
10 emulator_variable_precisions(9) = 10  ! Variable:      bp_r2d Routine:      ad_balance_tile Module:      ad_balance
11 emulator_variable_precisions(30) = 10 ! Variable:      dTdz Routine:      ad_balance_tile Module:      ad_balance
12 emulator_variable_precisions(31) = 10 ! Variable:      dSdz Routine:      ad_balance_tile Module:      ad_balance
13 emulator_variable_precisions(32) = 10 ! Variable:      ad_gradP Routine:      ad_balance_tile Module:      ad_balance
14 emulator_variable_precisions(33) = 10 ! Variable:      ad_phi Routine:      ad_balance_tile Module:      ad_balance
15 emulator_variable_precisions(34) = 10 ! Variable:      ad_gradPx Routine:      ad_balance_tile Module:      ad_balance
16 emulator_variable_precisions(35) = 10 ! Variable:      ad_gradPy Routine:      ad_balance_tile Module:      ad_balance
17 emulator_variable_precisions(36) = 10 ! Variable:      ad_A Routine:      ad_bc_r2d_tile Module:      ad_bc_2d
18 emulator_variable_precisions(37) = 10 ! Variable:      ad_A Routine:      ad_bc_u2d_tile Module:      ad_bc_2d
19 emulator_variable_precisions(38) = 10 ! Variable:      ad_A Routine:      ad_bc_v2d_tile Module:      ad_bc_2d
20 emulator_variable_precisions(39) = 10 ! Variable:      ad_A Routine:      ad_dabc_r2d_tile Module:      ad_bc_2d
21 emulator_variable_precisions(40) = 10 ! Variable:      ad_A Routine:      ad_dabc_u2d_tile Module:      ad_bc_2d
22 emulator_variable_precisions(41) = 10 ! Variable:      ad_A Routine:      ad_dabc_v2d_tile Module:      ad_bc_2d
23 emulator_variable_precisions(42) = 10 ! Variable:      ad_A Routine:      ad_bc_r2d_bry_tile Module:      ad_bc_bry2d
24 emulator_variable_precisions(43) = 10 ! Variable:      ad_A Routine:      ad_bc_u2d_bry_tile Module:      ad_bc_bry2d
25 emulator_variable_precisions(44) = 10 ! Variable:      ad_A Routine:      ad_bc_v2d_bry_tile Module:      ad_bc_bry2d
26 emulator_variable_precisions(45) = 10 ! Variable:      ad_A Routine:      ad_conv_r2d_bry_tile Module:      ad_conv_bry2d
27 emulator_variable_precisions(46) = 10 ! Variable:      ad_Awrk Routine:      ad_conv_r2d_bry_tile Module:      ad_conv_bry2d
28 emulator_variable_precisions(47) = 10 ! Variable:      ad_FE Routine:      ad_conv_r2d_bry_tile Module:      ad_conv_bry2d
29 emulator_variable_precisions(48) = 10 ! Variable:      ad_FX Routine:      ad_conv_r2d_bry_tile Module:      ad_conv_bry2d
30 emulator_variable_precisions(49) = 10 ! Variable:      Hfac Routine:      ad_conv_r2d_bry_tile Module:      ad_conv_bry2d

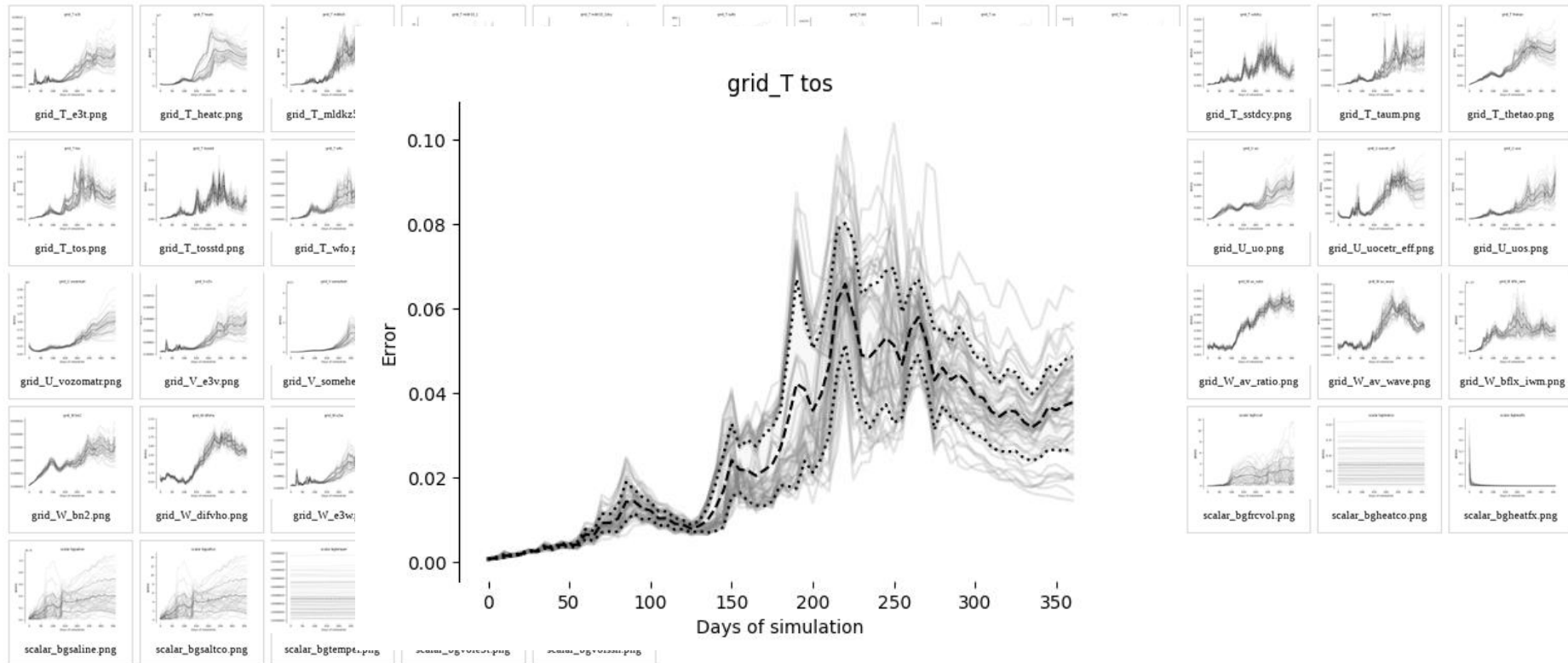
```

Precision Analysis

- How can we find which variables need to be kept in double-precision to maintain the accuracy of the results?:
 - **How can we measure the effect of reducing the precision of an arbitrary set of variables?**
 - How can we verify the accuracy of the results?

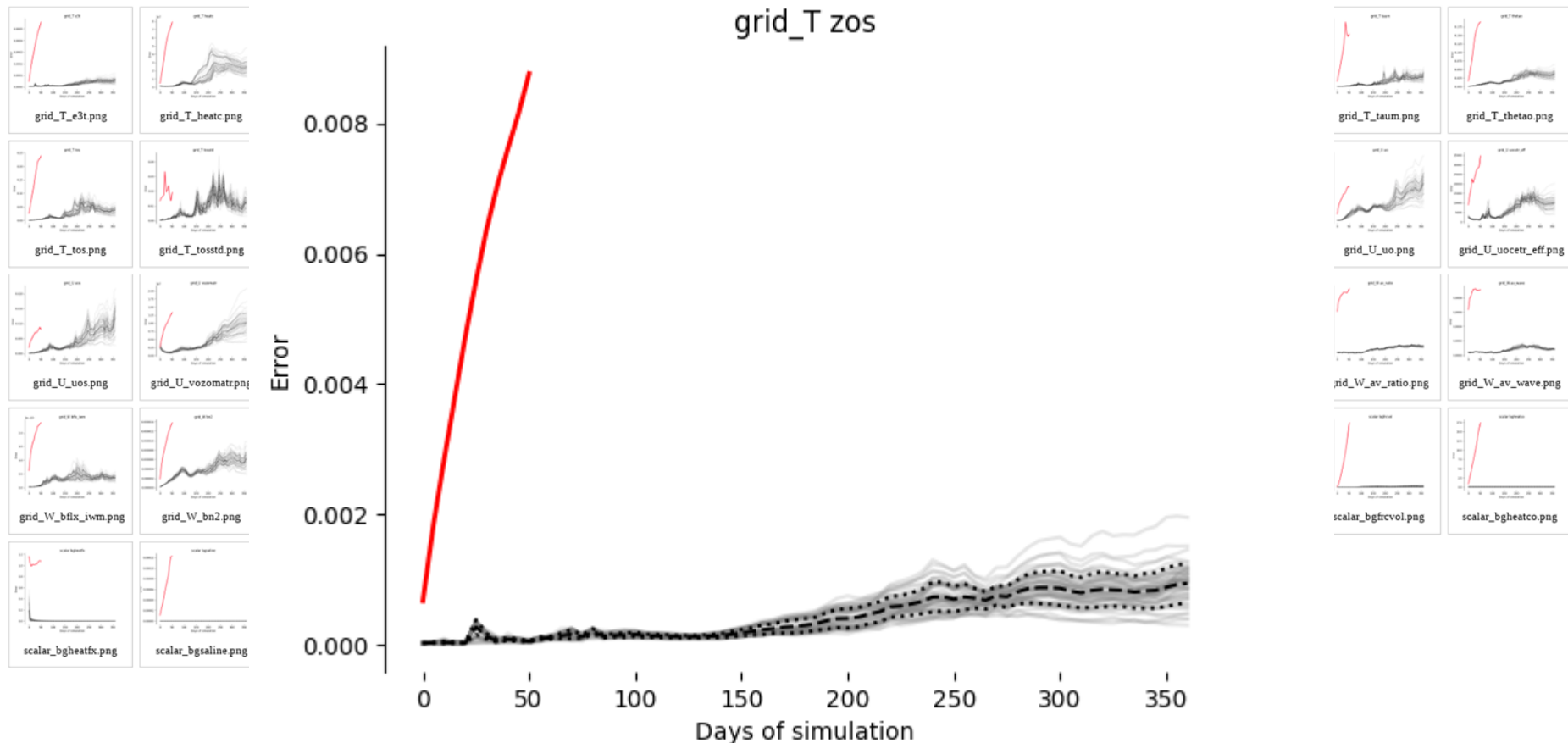
Verifying NEMO

- Initial conditions perturbed with white noise in the 3D temperature field.
- Evaluating 53 output variables.



Verifying NEMO

- Example: Everything in single precision:

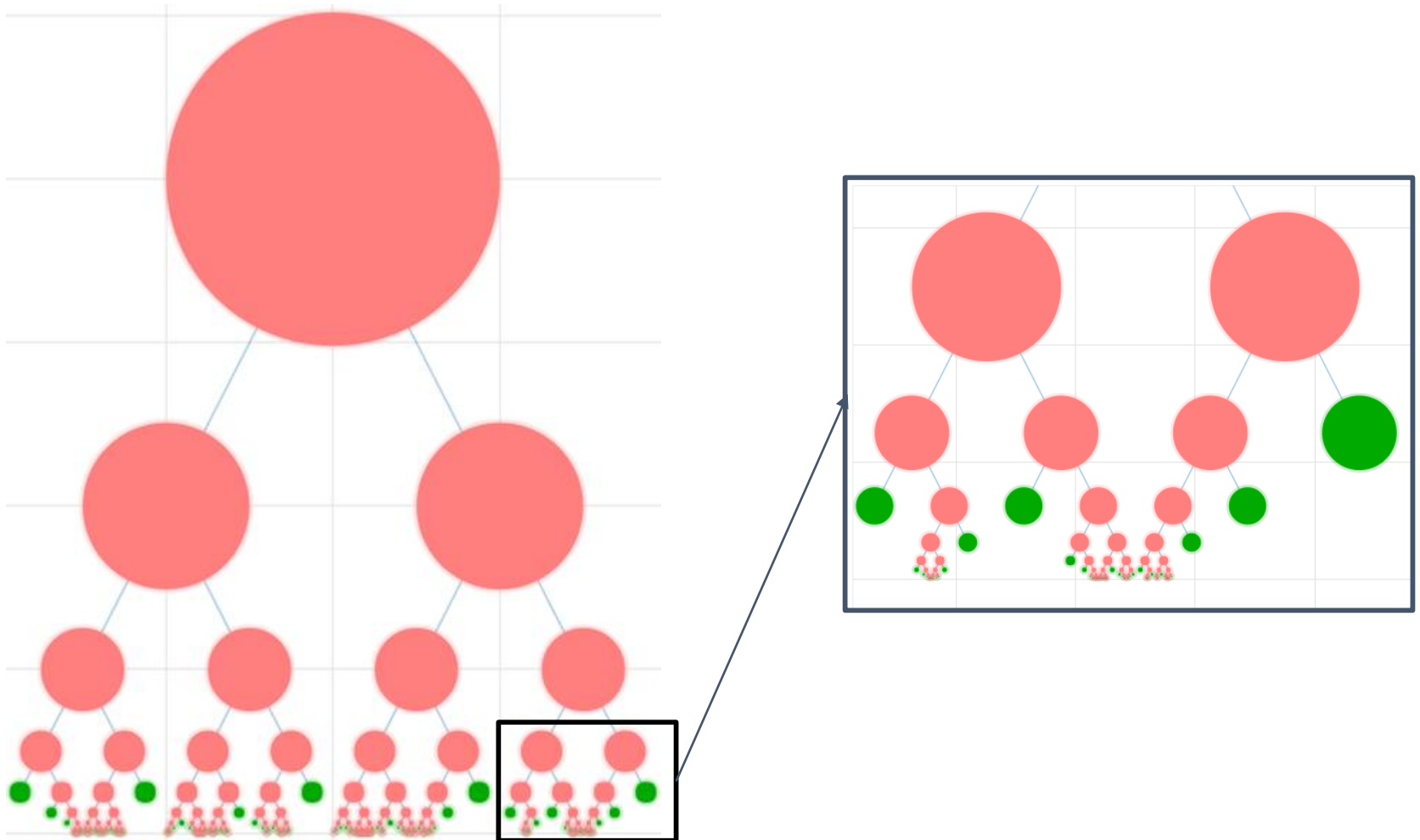


Precision Analysis

- How can we find which variables need to be kept in double-precision to maintain the accuracy of the results?:
 - How can we measure the effect of reducing the precision of an arbitrary set of variables?
 - How can we verify the accuracy of the results?

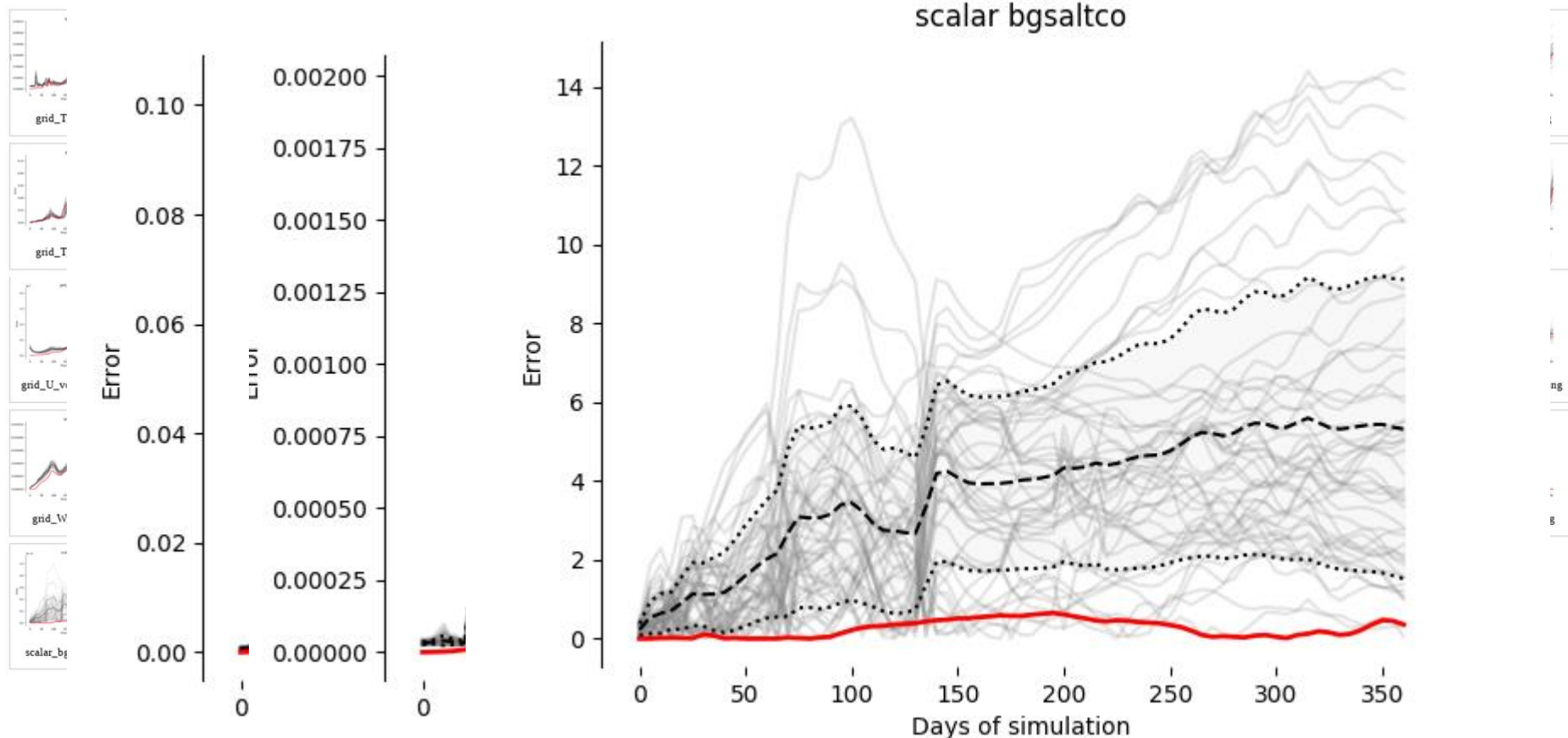
And now what?

Precision Analysis: NEMO



Results (1)

- Using this verification test to run the analysis algorithm on a small part of the code:

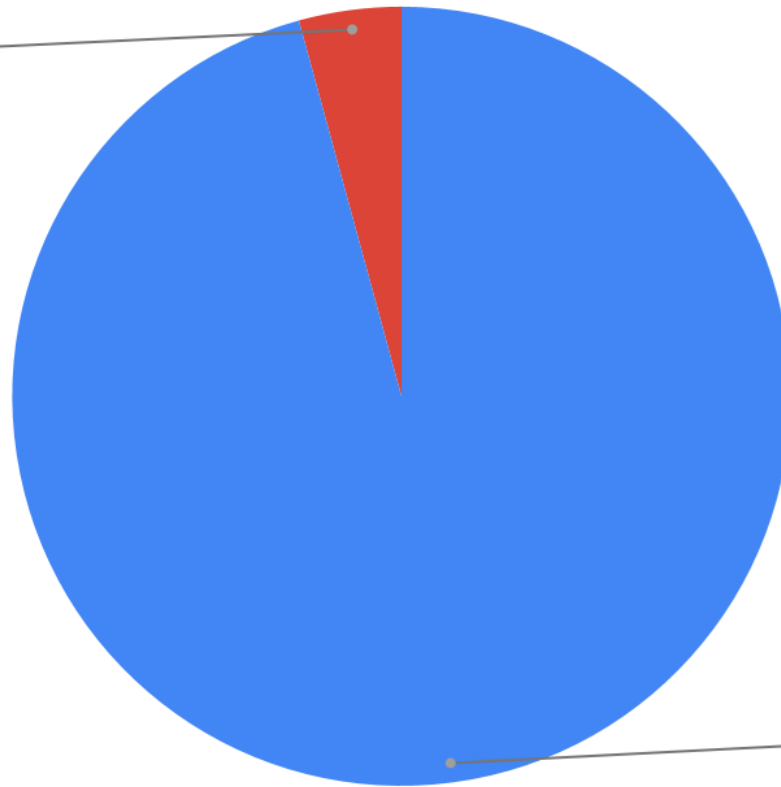


Results(2)

- Using the old verification test on the full model:

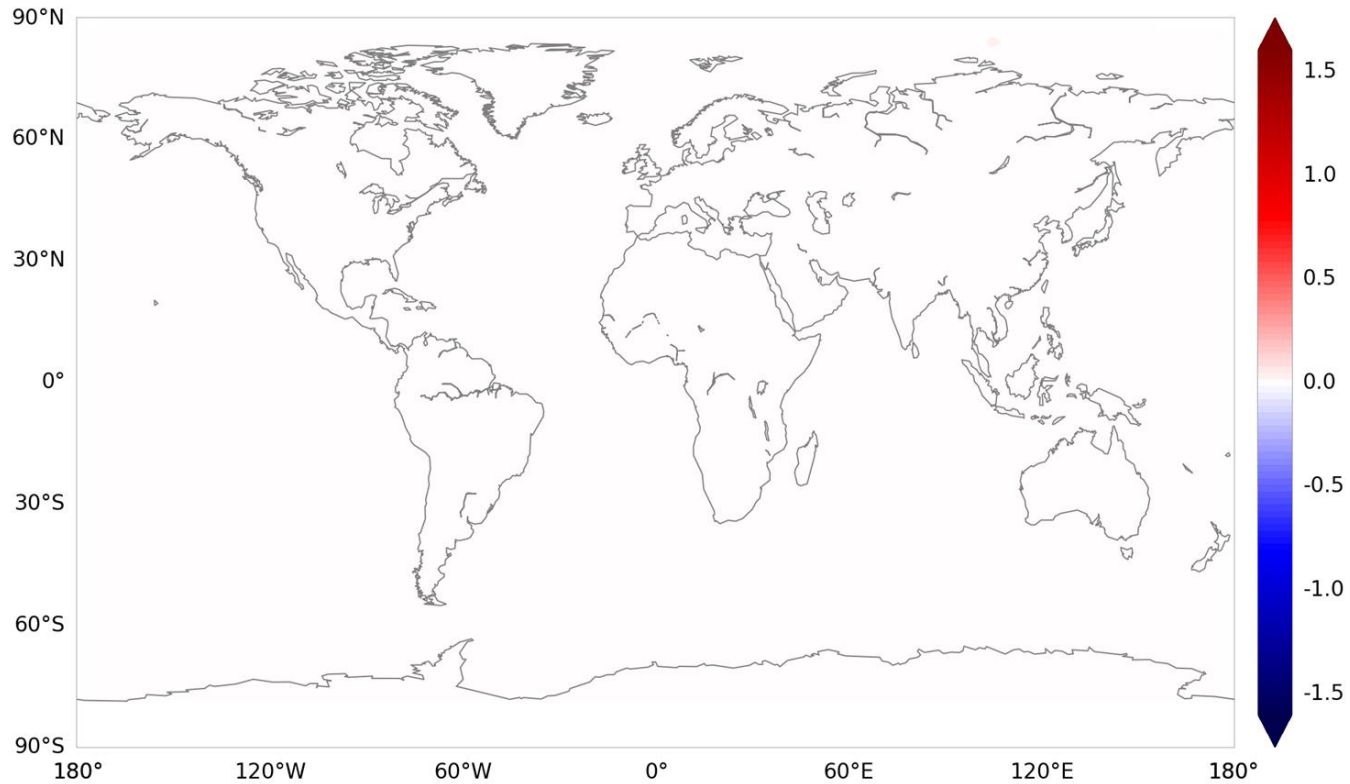
Results

Double Precision
4.2%



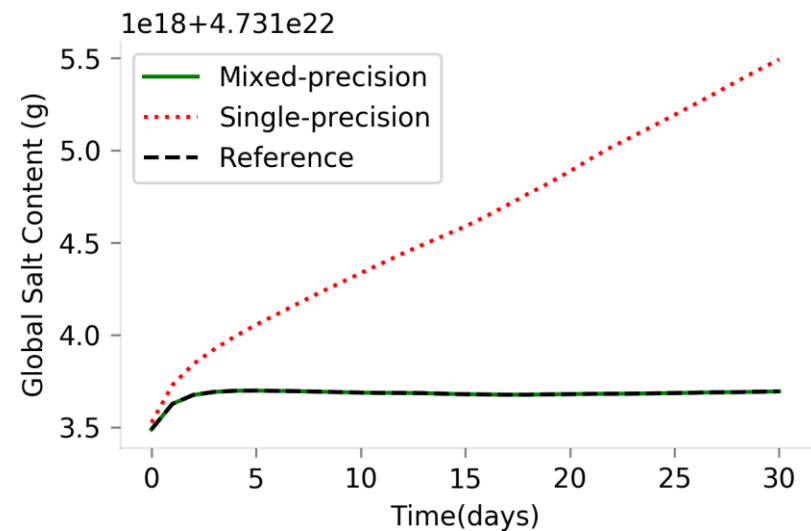
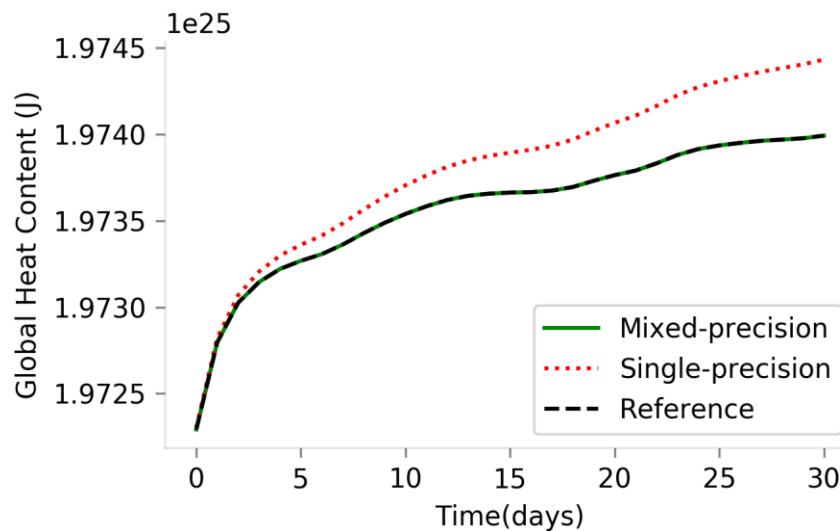
Single Precision
95.8%

But what about the outputs?



Difference between double- and mixed-precision monthly mean of sea-surface temperature for the first month of simulation.

But what about the outputs?



Conclusions

- Optimizing usage of numerical precision gives performance benefits.
- We can find which variables require double precision.
- There's a huge room for reducing the numerical precision in NEMO.

Ongoing work

- Performance in Mixed-precision?
- Transferability between different cases?
- Other ways of verifying the results?
- Reducing even more the precision for future architectures?



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación



**EXCELENCIA
SEVERO
OCHOA**



esiwace

CENTRE OF EXCELLENCE IN SIMULATION OF WEATHER
AND CLIMATE IN EUROPE

Thank you!

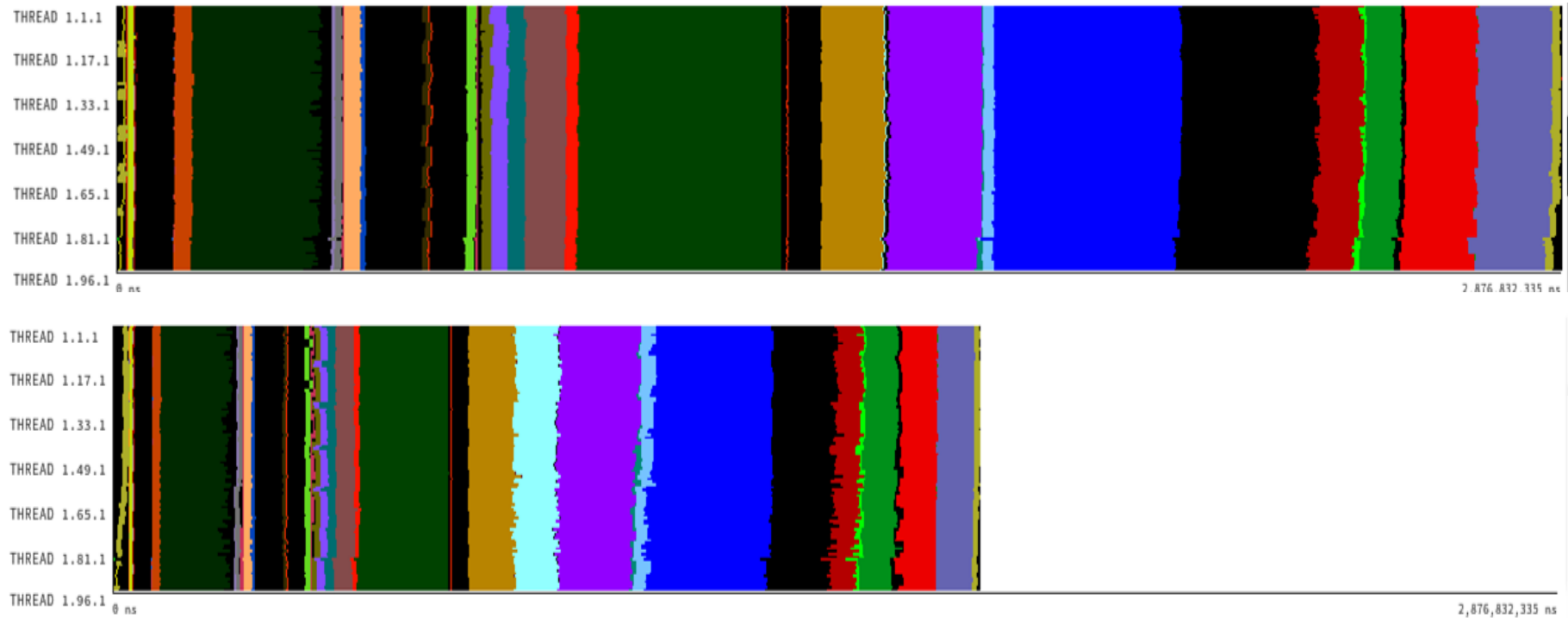
The research leading to these results has received funding from the EU H2020 Framework Programme under grant agreement H2020 GA 675191.

The content of this presentation reflects only the author's view. The European Commission is not responsible for any use that may be made of the information it contains.

oriol.tinto@bsc.es

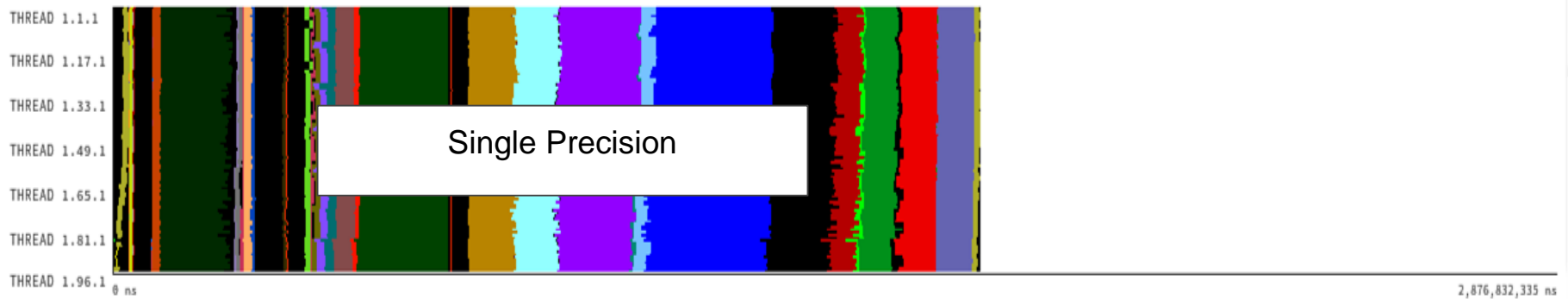
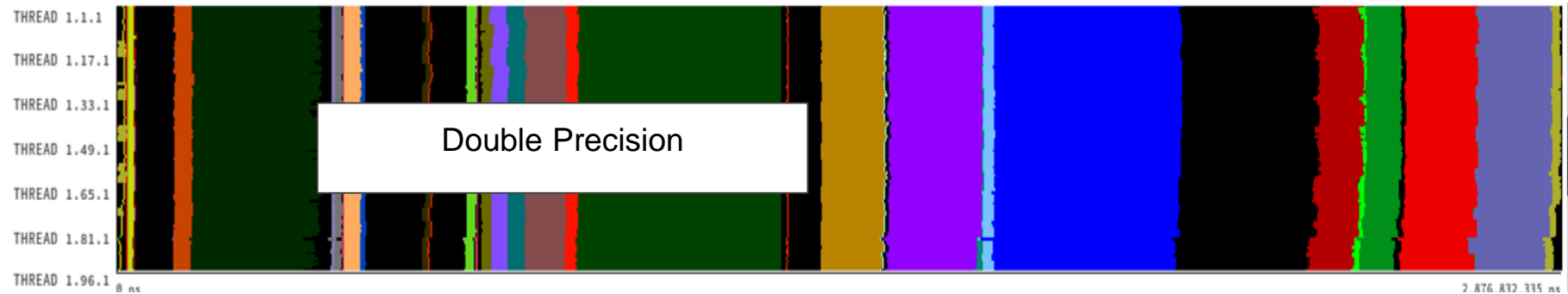
Bonus track:

How much a code can benefit from single precision?



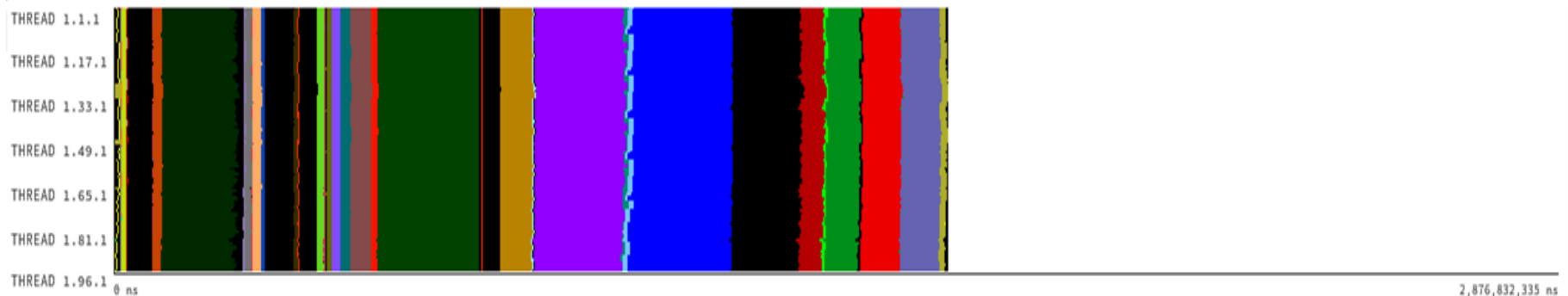
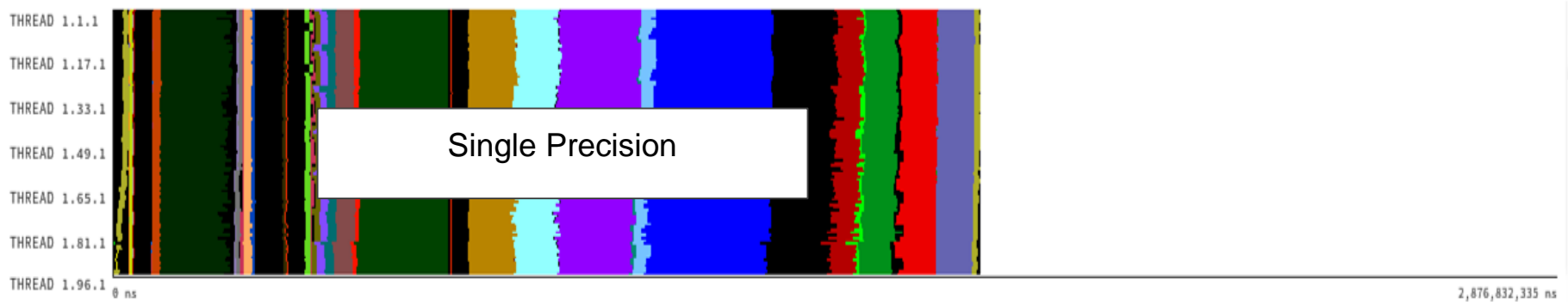
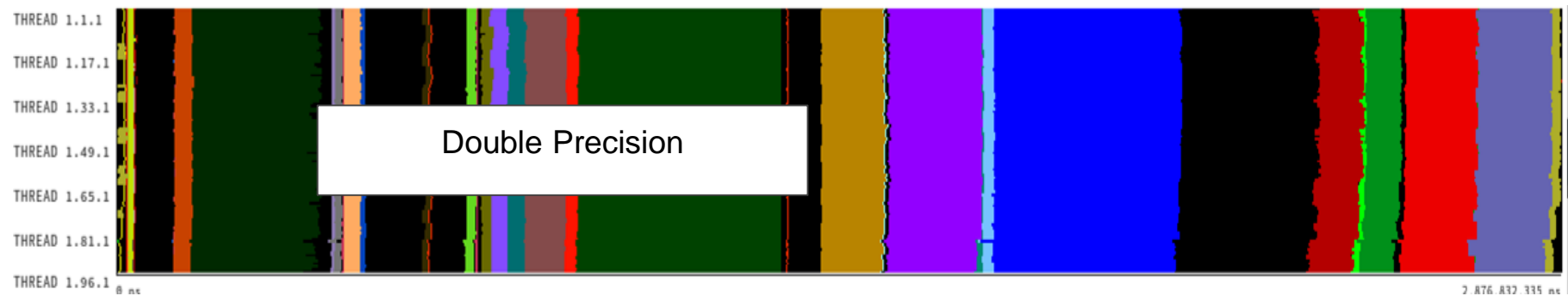
Bonus track:

How much a code can benefit from single precision?



Bonus track:

How much a code can benefit from single precision?



Bonus track:

How much a code can benefit from single precision?

